

Algorithmic stability for regression & classification

Rina Foygel Barber (joint with Jake Soloff & Rebecca Willett)

<http://rinafb.github.io/>

Collaborators



Jake Soloff



Rebecca Willett

Background on algorithmic stability

Background: algorithmic stability

Supervised learning setting

Data $(X_1, Y_1), \dots, (X_n, Y_n) \in \mathbb{R}^d \times \mathbb{R}$ $\xrightarrow{\text{algorithm } \mathcal{A}}$ Fitted model \hat{f}

Classical results are often of the form:

strong assumptions \implies strong guarantees
e.g., parametric model /
smoothness assumptions / etc e.g., $\hat{f} \approx$ true model,
with high probability

Background: algorithmic stability

More recently, assumption-lean / distribution-free results:

weaker assumptions

e.g., i.i.d. data from
an arbitrary distribution

\implies

weaker guarantees

e.g., can provide a prediction
interval around \hat{f}

Background: algorithmic stability

More recently, assumption-lean / distribution-free results:

weaker assumptions \implies weaker guarantees
e.g., i.i.d. data from an arbitrary distribution e.g., can provide a prediction interval around \hat{f}

- Generally do not guarantee or require concentration / consistency:

$\hat{f} \approx \hat{f}'$ if we resample entire data set

- But, often need a milder condition—stability:

$\hat{f} \approx \hat{f}'$ if we resample small fraction of data set

Background: algorithmic stability

Definition: algorithmic stability

\mathcal{A} is (ϵ, δ) -stable if, for any dataset \mathcal{D} of size n and any x ,

$$\frac{1}{n} \sum_{i=1}^n \mathbb{1} \left\{ \left| \hat{f}(x) - \hat{f}_{-i}(x) \right| \geq \epsilon \right\} \leq \delta$$

where $\hat{f} = \mathcal{A}(\mathcal{D})$ and $\hat{f}_{-i} = \mathcal{A}(\mathcal{D}_{-i})$

Intuition: only a small fraction of data points can be highly influential

Background: algorithmic stability

Definition: algorithmic stability

\mathcal{A} is (ϵ, δ) -stable if, for any dataset \mathcal{D} of size n and any x ,

$$\frac{1}{n} \sum_{i=1}^n \mathbb{P}_{\xi} \left\{ |\hat{f}(x) - \hat{f}_{-i}(x)| \geq \epsilon \right\} \leq \delta$$

where $\hat{f} = \mathcal{A}(\mathcal{D}; \xi)$ and $\hat{f}_{-i} = \mathcal{A}(\mathcal{D}_{-i}; \xi)$

random seed $\xi \sim \text{Unif}[0, 1]$

Intuition: only a small fraction of data points can be highly influential

Background: algorithmic stability

For statistical settings, a weaker definition is sufficient:

Definition: algorithmic stability with respect to P

For $\mathcal{D} \sim P^n$ and an independent test point $X \sim P_X$,

$$\frac{1}{n} \sum_{i=1}^n \mathbb{P} \left\{ |\hat{f}(X) - \hat{f}_{-i}(X)| \geq \epsilon \right\} \leq \delta$$

[Background: Kearns & Ron 1999; Bousquet & Elisseeff 2002; Elisseeff et al. 2005]

Background: algorithmic stability

Stability has implications for:

- Generalization [Bousquet & Elisseeff 2002; Elisseeff et al. 2005]
- Learnability [Shalev-Shwartz et al. 2010]
- Reproducibility [Yu 2013; Murdoch et al. 2019]
- Predictive inference with jackknife/cross-validation
[Steinberger & Leeb 2018; B., Candès, Ramdas, Tibshirani 2021]
- Asymptotics of cross-validation
[Kumar et al. 2013; Austern & Zhou 2020; Bayle et al. 2020]

Can algorithmic stability be assumed?

Some algorithms are stable by construction, e.g., k -NN, ridge regression

[Rogers & Wagner 1978; Bousquet & Elisseeff 2002]

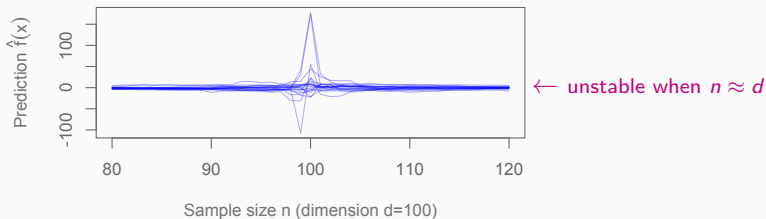
Can algorithmic stability be assumed?

Some algorithms are stable by construction, e.g., k -NN, ridge regression

[Rogers & Wagner 1978; Bousquet & Elisseeff 2002]

But, even simple algorithms can be unstable

- Least squares: $\hat{f}(x) = x^T \hat{\theta}$ where $\hat{\theta} = (\mathbf{X}^T \mathbf{X})^\dagger \mathbf{X}^T \mathbf{Y}$



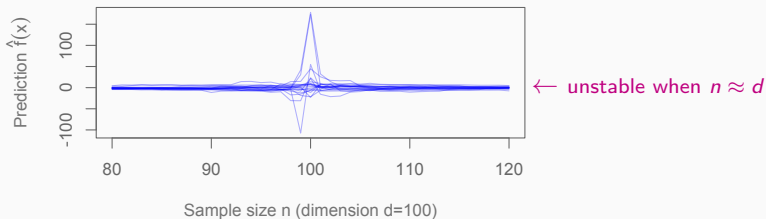
Can algorithmic stability be assumed?

Some algorithms are stable by construction, e.g., k -NN, ridge regression

[Rogers & Wagner 1978; Bousquet & Elisseeff 2002]

But, even simple algorithms can be unstable

- Least squares: $\hat{f}(x) = x^T \hat{\theta}$ where $\hat{\theta} = (\mathbf{X}^T \mathbf{X})^\dagger \mathbf{X}^T \mathbf{Y}$



And, many modern algorithms are too complex for theoretical analysis

Can algorithmic stability be tested?

If stability is not known theoretically — can we certify it empirically?

The “black-box” setting:

We can only learn how \mathcal{A} works by running it on data, e.g.,

- subsamples / bootstrapped samples of the real data
- perturbations of the real data
- simulate data from a model fitted to real data

Can algorithmic stability be tested?

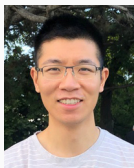
Hardness result for black-box tests:

To certify that \mathcal{A} is stable with respect to P at sample size n ,

- Need available sample size $\gg n$ drawn from P
- Or, need to perform exhaustive search over the data space



Byol Kim



Yuetian Luo

Kim & B. 2021, *Black-box tests for algorithmic stability*, arXiv:2111.15546

Luo & B. 2024, *Is algorithmic stability testable? A unified framework under computational constraints*, arXiv:2405.15107

Key question

Since algorithmic stability can't be certified empirically
(and we don't want to always use simple algorithms like k -NN)

Key question

Is there a way to convert any algorithm into a stable algorithm, while still retaining its good performance?

Key question

Our work uses *bagging* to address this question in a range of settings:

- Part 1: stability for regression (real-valued response)
- Part 2: stability in Hilbert & Banach spaces (e.g., functions)
- Part 3: stability for classification



Jake Soloff



Rebecca Willett

Bagging provides assumption-free stability, [arxiv:2301.12600](https://arxiv.org/abs/2301.12600)

Stability via resampling: statistical problems beyond the real line, [arxiv:2405.09511](https://arxiv.org/abs/2405.09511)

Building a stable classifier with the inflated argmax, [arxiv:2405.14064](https://arxiv.org/abs/2405.14064)


Part 1: stability for bagged algorithms in regression

Background on bagging

Empirically, bagging (& other ensembling procedures) have been observed to improve stability dramatically.

Bagging

- Sample subsets of data ("bags") r_b for $b = 1, \dots, B$
- Fit models $\hat{f}_b = \mathcal{A}(\mathcal{D}_{r_b}; \xi_b)$ where $\mathcal{D}_{r_b} = ((X_i, Y_i) : i \in r_b)$

 random seed

- $\tilde{\mathcal{A}}_B$ returns aggregated model \hat{f} :

$$\hat{f}(x) := \frac{1}{B} \sum_b \hat{f}_b(x)$$

Background on bagging

The most common options:

- Bootstrapping (“bagging”) = subsets r_b of size m , sampled uniformly with replacement

$$p = \mathbb{P}\{i \in r_b\} = 1 - (1 - 1/n)^m$$

- Subbagging = subsets r_b of size $m < n$, sampled uniformly without replacement

$$p = \mathbb{P}\{i \in r_b\} = m/n$$

[Breiman 1996; Andonova et al. 2002]

Background on bagging

Bagging appears in:

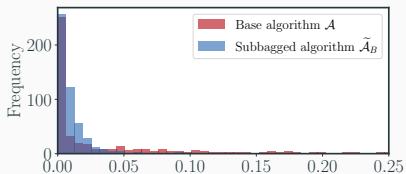
- Random forests [Breiman 2001]
- Variable selection in regression [Meinshausen & Bühlmann 2010]
- Classification in the presence of class imbalance [Li 2007]
- Robust Bayesian inference (BayesBag) [Huggins & Miller 2023]
- & many more

Many results on theoretical properties—

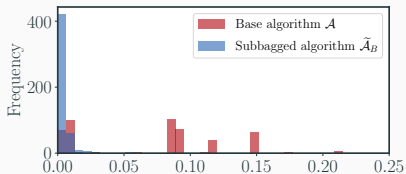
Bagging induces smoothness, reduces variance, creates robustness

[Bühlmann & Yu 2002; Grandvalet 2004; Friedman & Hall 2000; Elisseeff et al. 2005]

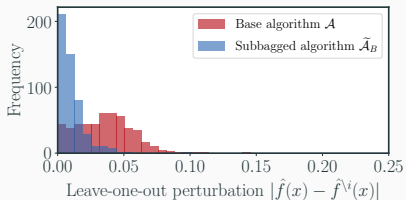
Bagging helps stability



← logistic regression



← regression trees



← neural networks

Main result: stability guarantee

Recall definition of (ϵ, δ) -stability:

$$\frac{1}{n} \sum_{i=1}^n \mathbb{P}_{\xi} \left\{ |\hat{f}(x) - \hat{f}_{-i}(x)| \geq \epsilon \right\} \leq \delta$$

Main result: stability guarantee

Recall definition of (ϵ, δ) -stability:

$$\frac{1}{n} \sum_{i=1}^n \mathbb{P}_{\xi} \left\{ |\hat{f}(x) - \hat{f}_{-i}(x)| \geq \epsilon \right\} \leq \delta$$

Theorem

Let \mathcal{A} be any base algorithm that returns predictions in $[0, 1]$.
Then $\tilde{\mathcal{A}}_B$ satisfies (ϵ, δ) -stability as long as

$$\epsilon^2 \delta \geq \frac{1}{4(n-1)} \cdot \frac{p}{1-p} + \mathcal{O}(B^{-1})$$

Main result: stability guarantee

- Interpretation: the bagged version of *any* base alg. is always stable:
 - No assumptions on data—holds for all \mathcal{D} and all x
 - No assumptions on \mathcal{A} aside from boundedness
- Can also extend to algorithms with unbounded output, via either clipping predictions, or adapting to the range of \hat{f} 's

Defining stability: the mean or the tail?

Tail stability

\mathcal{A} is (ϵ, δ) -stable if, for any dataset \mathcal{D} of size n and any x ,

$$\frac{1}{n} \sum_{i=1}^n \mathbb{P}_{\xi} \left\{ |\hat{f}(x) - \hat{f}_{-i}(x)| \geq \epsilon \right\} \leq \delta$$

Mean-square stability

\mathcal{A} is β^2 -stable if, for any dataset \mathcal{D} of size n and any x ,

$$\frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\xi} \left[(\hat{f}(x) - \hat{f}_{-i}(x))^2 \right] \leq \beta^2$$

Defining stability: the mean or the tail?

$\underbrace{\beta^2\text{-stability}}_{\text{mean-square stability}} \quad \text{vs} \quad \underbrace{(\epsilon, \delta)\text{-stability}}_{\text{tail stability}}$

- $\beta^2\text{-stability} \implies (\epsilon, \delta)\text{-stability}$ for any $\epsilon^2\delta \geq \beta^2$
- $(\epsilon, \delta)\text{-stability, \& outputs} \in [0, 1] \implies \beta^2\text{-stab.}$ for $\beta^2 = \epsilon^2 + \delta$

Main result: stability guarantee

Theorem

Let \mathcal{A} be any base algorithm that returns predictions in $[0, 1]$.
Then $\tilde{\mathcal{A}}_B$ satisfies β^2 -stability for

$$\beta^2 = \frac{1}{4(n-1)} \cdot \frac{\rho}{1-\rho} + \mathcal{O}(B^{-1})$$

- Directly implies the (ϵ, δ) -stability result

Main result: stability guarantee

Regimes for subbagging (with $B \rightarrow \infty$):

- **Proportional sampling:**

$$m = \mathcal{O}(n) \implies \beta^2 \propto n^{-1}$$

Main result: stability guarantee

Regimes for subbagging (with $B \rightarrow \infty$):

- **Proportional sampling:**

$$m = \mathcal{O}(n) \implies \beta^2 \propto n^{-1}$$

- **Massive subsampling:**

$$m = \mathcal{O}(n^a) \implies \beta^2 \propto n^{-(2-a)}$$


[See also Poggio et al. 2002; Elisseff et al. 2005; Chen et al. 2022]

Main result: stability guarantee

Regimes for subbagging (with $B \rightarrow \infty$):

- **Proportional sampling:**

$$m = \mathcal{O}(n) \implies \beta^2 \propto n^{-1}$$

- **Massive subsampling:**

$$m = \mathcal{O}(n^a) \implies \beta^2 \propto n^{-(2-a)}$$

0 < a < 1

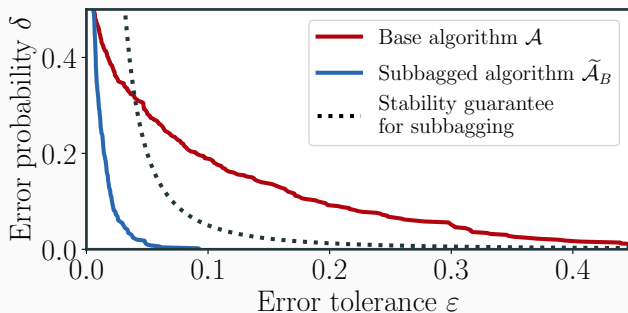
[See also Poggio et al. 2002; Elisseff et al. 2005; Chen et al. 2022]

- **Minimal subsampling:**

$$m = n - \mathcal{O}(n^a) \implies \beta^2 \propto n^{-a}$$

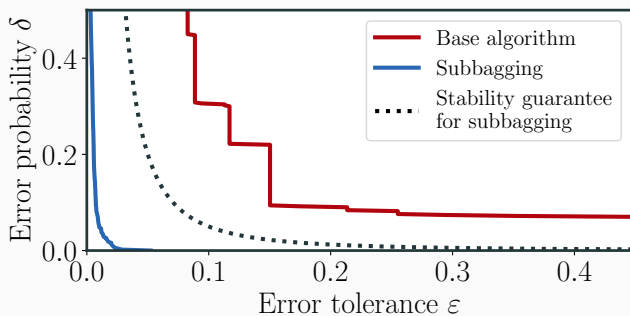
Experiments

Logistic regression ($d = 200$, $n = 500$, $m = n/2$, $B = 10000$)



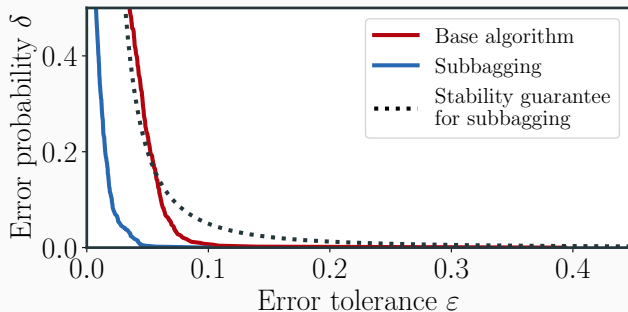
Experiments

Regression trees ($d = 40$, $n = 500$, $m = n/2$, $B = 10000$)



Experiments

Neural networks ($d = 200$, $n = 500$, $m = n/2$, $B = 10000$)



Is the guarantee tight?

Could the bound be improved?

Theorem: a matching bound for subbagging

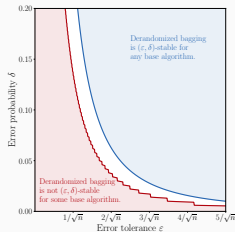
Fix n, ρ with $\min\{\rho, 1 - \rho\} \gg \frac{1}{n}$.

There exists an \mathcal{A} with output in $[0, 1]$, such that if

$$\beta^2 < \frac{1}{c(n-1)} \cdot \frac{\rho}{1-\rho}$$

then $\tilde{\mathcal{A}}_\infty$ is *not* β^2 -stable.

Illustration for $n = 500$, $m = 250$:



Is the guarantee tight?

Could we use a stricter definition of stability?

Worst-case stability

\mathcal{A} is β -worst-case-stable if, for any dataset \mathcal{D} of size n and any x ,

$$\max_i |\hat{f}(x) - \hat{f}_{-i}(x)| \leq \beta$$

max instead of average

Is the guarantee tight?

Could we use a stricter definition of stability?

Worst-case stability

\mathcal{A} is β -worst-case-stable if, for any dataset \mathcal{D} of size n and any x ,

$$\max_i |\hat{f}(x) - \hat{f}_{-i}(x)| \leq \beta$$

max instead of average 

- For any \mathcal{A} with output in $[0, 1]$, $\tilde{\mathcal{A}}_\infty$ is β -worst-case-stable for $\beta = p$ [Poggio et al. 2002; Elisseff et al. 2005; Chen et al. 2022]
- For any $\beta < p$, can construct \mathcal{A} s.t. $\tilde{\mathcal{A}}_\infty$ is not β -worst-case-stable

Part 2: bagged algorithms beyond the real line

Settings & applications

The stability result proved so far....

- Stated for prediction: stability of $\mathcal{D} \mapsto [\mathcal{A}(\mathcal{D})](x) \in \mathbb{R}$
- Result holds for *any* real-valued output, i.e., $\mathcal{D} \mapsto \mathcal{A}(\mathcal{D}) \in \mathbb{R}$
(doesn't have to be “a prediction”)

Settings & applications

The stability result proved so far....

- Stated for prediction: stability of $\mathcal{D} \mapsto [\mathcal{A}(\mathcal{D})](x) \in \mathbb{R}$
- Result holds for *any* real-valued output, i.e., $\mathcal{D} \mapsto \mathcal{A}(\mathcal{D}) \in \mathbb{R}$
(doesn't have to be “a prediction”)

Many data analysis algorithms return output that is more complex:

- Estimate a discrete distribution $\rightsquigarrow \mathcal{A}(\mathcal{D}) \in \Delta_{L-1} \subseteq \mathbb{R}^L$ (simplex)
- Estimate a density on $\mathbb{R} \rightsquigarrow \mathcal{A}(\mathcal{D}) \in L_1(\mathbb{R})$
- Fit a parametric model $\rightsquigarrow \mathcal{A}(\mathcal{D}) \in \Theta \subseteq \mathbb{R}^d$

Stability guarantee for a Hilbert space

Suppose \mathcal{A} returns outputs lying in a Hilbert space:

$$\mathcal{A} : \bigcup_{n \geq 0} \mathcal{Z}^n \longrightarrow \mathcal{W} \subseteq \mathcal{H}$$

Hilbert space with norm $\|\cdot\|_{\mathcal{H}}$

Stability guarantee for a Hilbert space

Suppose \mathcal{A} returns outputs lying in a Hilbert space:

$$\mathcal{A}: \bigcup_{n \geq 0} \mathcal{Z}^n \longrightarrow \mathcal{W} \subseteq \mathcal{H}$$

↖
Hilbert space with norm $\|\cdot\|_{\mathcal{H}}$

Assume \mathcal{W} is a bounded subset:

$$\text{rad}_{\mathcal{H}}(\mathcal{W}) = \inf_{w \in \mathcal{W}} \sup_{w' \in \mathcal{W}} \|w - w'\|_{\mathcal{H}}$$

- Example: $\mathcal{W} = \Delta_{L-1} \subseteq \mathbb{R}^L \rightsquigarrow \text{rad}_{\mathcal{H}}(\mathcal{W})^2 = 1 - \frac{1}{L}$

Stability guarantee for a Hilbert space

Definition

\mathcal{A} is β^2 -stable with respect to $\|\cdot\|_{\mathcal{H}}$ if, for any dataset \mathcal{D} of size n ,

$$\frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\xi} [\|\mathcal{A}(\mathcal{D}; \xi) - \mathcal{A}(\mathcal{D}_{-i}; \xi)\|_{\mathcal{H}}^2] \leq \beta^2$$

Theorem

Let \mathcal{A} be any base algorithm that returns predictions in $\mathcal{W} \subseteq \mathcal{H}$.

Then $\tilde{\mathcal{A}}_B$ satisfies β^2 -stability with respect to $\|\cdot\|_{\mathcal{H}}$ for

$$\beta^2 = \frac{\text{rad}_{\mathcal{H}}(\mathcal{W})^2}{n-1} \cdot \frac{p}{1-p} + \mathcal{O}(B^{-1}).$$

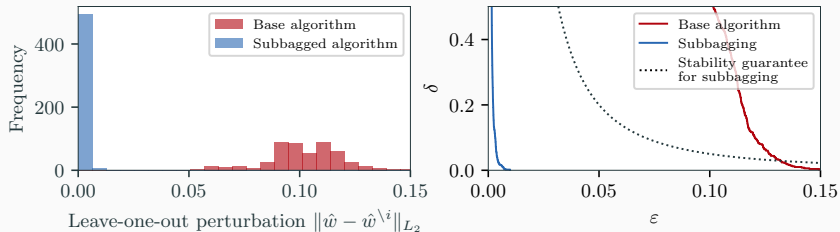
Stability guarantee for a Hilbert space

- No price for high dimensionality
- Earlier result for $[0, 1] \subseteq \mathbb{R}$ is a special case, with $\text{rad}_{\mathcal{H}}(\mathcal{W}) = 1/2$
- Can extend to Banach space (e.g., d_{TV} distance) $\rightsquigarrow \log(\text{dim})$ terms

Experiments

Estimating a regression function in L_2 (base algorithm = regression tree)

Setting: $d = 40$, $n = 500$, $m = n/2$, $B = 10000$



Part 3: a stability framework for classification

Can classification be stable?

Returning to the supervised learning setting...

- Training data $\mathcal{D} = ((X_1, Y_1), \dots, (X_n, Y_n))$
- Test point x
- $\mathcal{Y} = [L]$ = finite set of labels

If we use \mathcal{D} to train a classifier (a map $x \mapsto \hat{y} \in [L]$)...

- Instability must occur b/c $\hat{y} \neq \hat{y}_{-i}$ is likely for ambiguous examples

Can classification be stable?

Returning to the supervised learning setting...

- Training data $\mathcal{D} = ((X_1, Y_1), \dots, (X_n, Y_n))$
- Test point x
- $\mathcal{Y} = [L] =$ finite set of labels

If we use \mathcal{D} to train a classifier (a map $x \mapsto \hat{y} \in [L]$)...

- Instability must occur b/c $\hat{y} \neq \hat{y}_{-i}$ is likely for ambiguous examples
- A common approach: return a set of candidate labels, $x \mapsto \hat{S} \subseteq [L]$
[Grycko 1993; del Coz et al. 2009; Sadinle et al. 2019; Chzhen et al. 2021]

Can classification be stable?

Selection stability

A classification algorithm satisfies selection stability at level δ if

$$\frac{1}{n} \sum_{i=1}^n \mathbb{1} \left\{ \widehat{S} \cap \widehat{S}_{-i} = \emptyset \right\} \leq \delta$$

where $\widehat{S} = [\mathcal{A}(\mathcal{D})](x)$ and $\widehat{S}_{-i} = [\mathcal{A}(\mathcal{D}_{-i})](x)$

Interpretation: if $\widehat{S} \cap \widehat{S}_{-i} \neq \emptyset$, then it's possible to have

$$y_* \in \widehat{S} \text{ and } y_* \in \widehat{S}_{-i}$$

where y_* is the true label

Can classification be stable?

Selection stability

A classification algorithm satisfies selection stability at level δ if

$$\frac{1}{n} \sum_{i=1}^n \mathbb{P}_{\xi} \left\{ \hat{S} \cap \hat{S}_{-i} = \emptyset \right\} \leq \delta$$

where $\hat{S} = [\mathcal{A}(\mathcal{D}; \xi)](x)$ and $\hat{S}_{-i} = [\mathcal{A}(\mathcal{D}_{-i}; \xi)](x)$

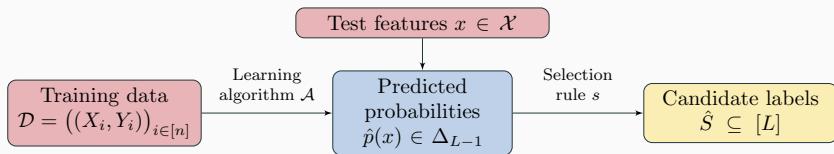
Interpretation: if $\hat{S} \cap \hat{S}_{-i} \neq \emptyset$, then it's possible to have

$$y_* \in \hat{S} \text{ and } y_* \in \hat{S}_{-i}$$

where y_* is the true label

A two-stage framework

Building a classification algorithm:



A two-stage framework

Definition: ϵ -compatible selection rule

If $w, w' \in \Delta_{L-1}$ with $\|w - w'\|_2 < \epsilon$, then $s(w) \cap s(w') \neq \emptyset$

Stability in the two-stage framework

Any (ϵ, δ) -stable \mathcal{A} + any ϵ -compatible $s \Rightarrow \delta$ -selection-stable $s \circ \mathcal{A}$

A two-stage framework

Definition: ϵ -compatible selection rule

If $w, w' \in \Delta_{L-1}$ with $\|w - w'\|_2 < \epsilon$, then $s(w) \cap s(w') \neq \emptyset$

Stability in the two-stage framework

Any (ϵ, δ) -stable \mathcal{A} + any ϵ -compatible $s \Rightarrow \delta$ -selection-stable $s \circ \mathcal{A}$

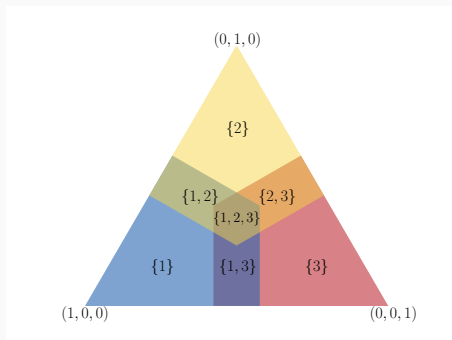
construct via
bagging any base alg.

how to define s ?

A baseline: the fixed-margin rule

A fixed-margin selection rule is ϵ -compatible:

$$s_{\text{margin}}^{\epsilon}(w) = \left\{ j \in [L] : w_j > \max_k w_k - \epsilon/\sqrt{2} \right\}$$



The inflated argmax

Is the fixed-margin rule optimal?

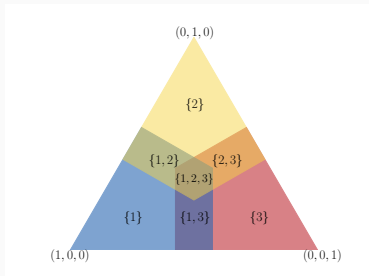
(Does it return small sets whenever possible?)

The inflated argmax

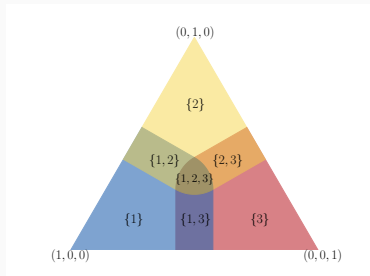
Is the fixed-margin rule optimal?

(Does it return small sets whenever possible?)

Comparison:



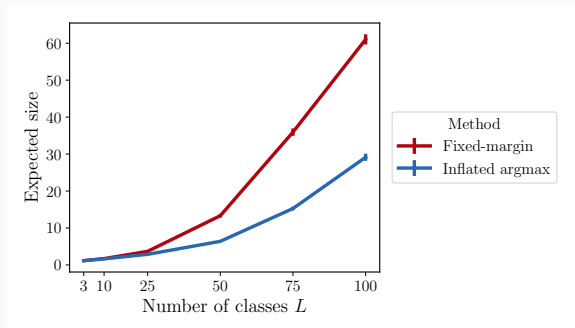
Fixed-margin rule



Inflated argmax

The inflated argmax

Another comparison:



The inflated argmax

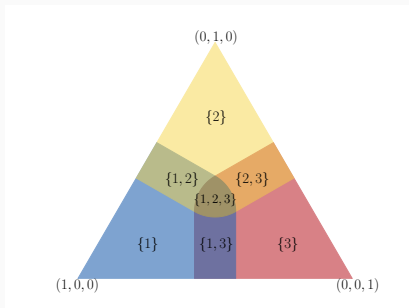
The inflated argmax selection rule

Define

$$\arg \max^\epsilon(w) = \{j \in [L] : \text{dist}_{\ell_2}(w, R_j) < \epsilon\},$$

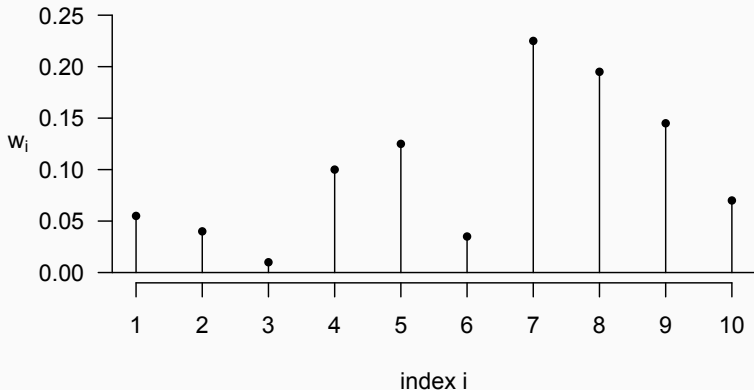
where

$$R_j = \left\{ w \in \mathbb{R}^L : w_j \geq \max_{k \neq j} w_k + \epsilon/\sqrt{2} \right\}$$



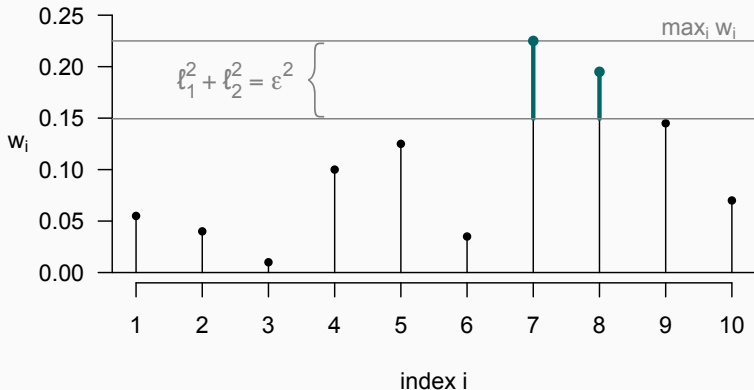
The inflated argmax

Computing the inflated argmax:



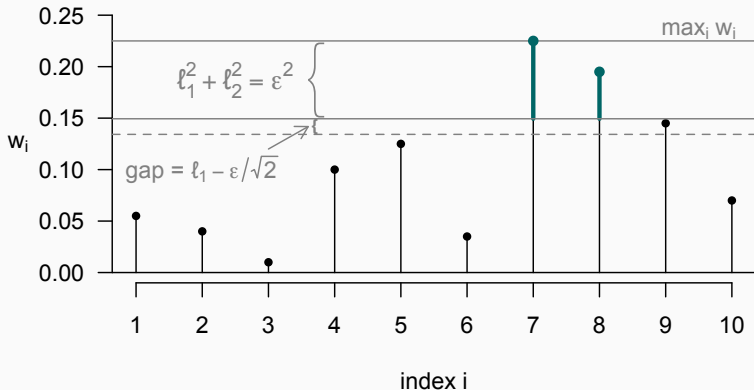
The inflated argmax

Computing the inflated argmax:



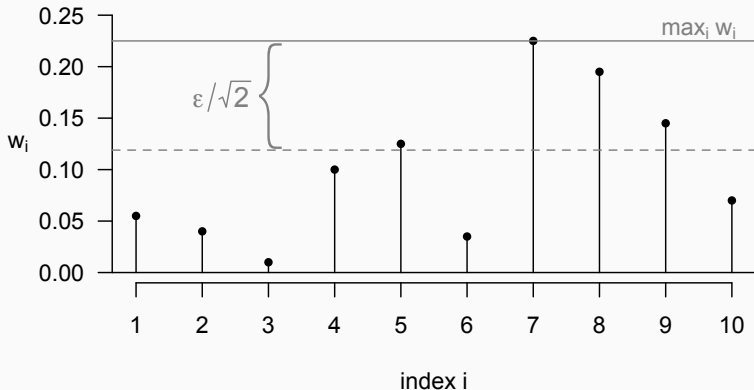
The inflated argmax

Computing the inflated argmax:



The inflated argmax

Compare to the fixed-margin selection rule:



Properties of the inflated argmax

Theorem (compatibility)

The inflated argmax is ϵ -compatible:

$$\|w - w'\|_2 < \epsilon \implies \arg \max^\epsilon(w) \cap \arg \max^\epsilon(w') \neq \emptyset$$

Interpretation: any bagged alg. + inflated argmax \rightsquigarrow selection stability

Properties of the inflated argmax

Theorem (compatibility)

The inflated argmax is ϵ -compatible:

$$\|w - w'\|_2 < \epsilon \implies \arg \max^\epsilon(w) \cap \arg \max^\epsilon(w') \neq \emptyset$$

Interpretation: any bagged alg. + inflated argmax \rightsquigarrow selection stability

Theorem (optimality)

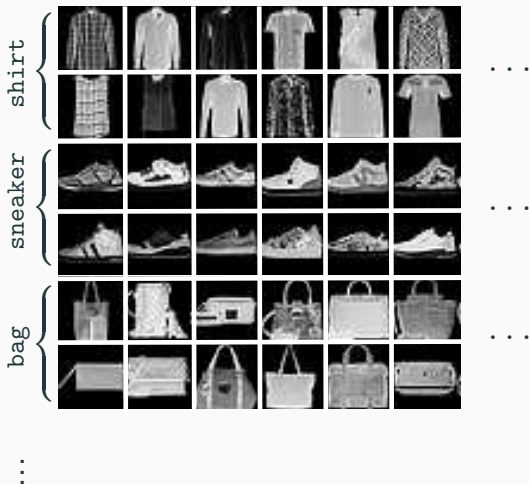
For any ϵ -compatible selection rule s ,
if s is symmetric and $s(w) \supseteq \arg \max(w)$ for all w ,

$$s(w) = \{j\} \implies \arg \max^\epsilon(w) = \{j\}$$

Interpretation: inflated argmax returns a singleton set as often as possible

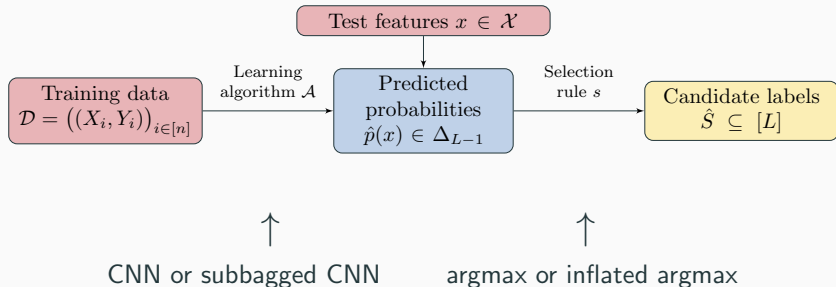
Experiments

Dataset: fashion MNIST ($L = 10$ classes) [figure & data from Xiao et al. 2017]



Experiments

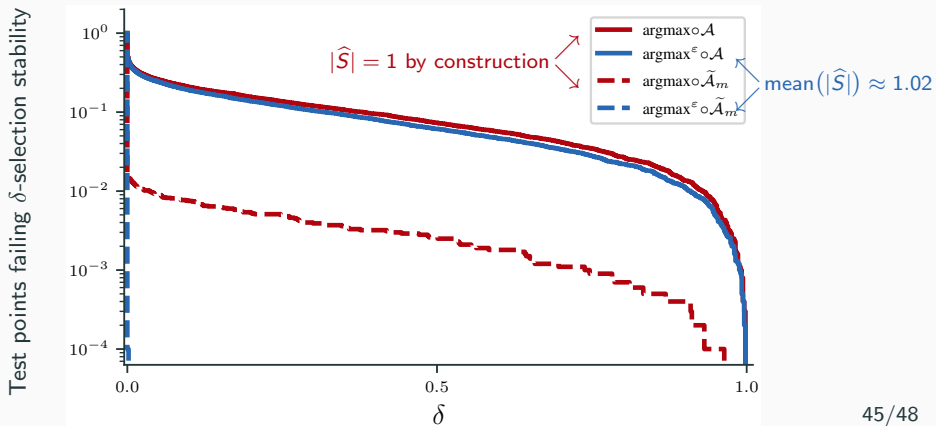
Methods tested:



Experiments

Setting: $n = 60000$, $m = n/2$, $B = 200$, inflation $\epsilon = 0.05$

Methods: {CNN or subbagged CNN} + {argmax or inflated argmax}



An extension: stability for ranking problems

Data from n users \rightsquigarrow scores $\hat{w} = (\hat{w}_1, \dots, \hat{w}_L)$ for L items.

- What are the top k items?
(return an inflated set of size $\geq k$)
- What is the full ranked list?
(return ≥ 1 rankings of the list)

	Godfather	The Matrix	Toy Story	...
User #1		5	4	...
User #2	3			...
User #3		1	3	...
\vdots	\vdots	\vdots	\vdots	



Ruiting Liang



Jake Soloff



Rebecca Willett

Summary & open questions

Summary & open questions

Our results:

- Classical bagging & subbagging can be applied to *any* algorithm \mathcal{A} to achieve an assumption-free stability guarantee
- Downstream, leads to generalization, predictive inference, etc
- Extends to outputs in a Hilbert or Banach space (e.g., output can be a vector, a function, a distribution,)
- Can combine with inflated argmax for stable classification

Summary & open questions

Open questions:

- How does bagging perform relative to other definitions of stability, & related properties — generalization, robustness, privacy, etc?
 - Guarantees for aggregation procedures aside from averaging?
 - Is the computational cost of bagging a necessary cost for stability?
-

Thank you!