

## 3.3 How to Solve Convex OPs and SVM

*Machine Learning 1: Foundations*

Marius Kloft (TUK)

# The Good News

# The Good News

Convex OPs are **easy to solve!**

# The Good News

Convex OPs are **easy to solve!**

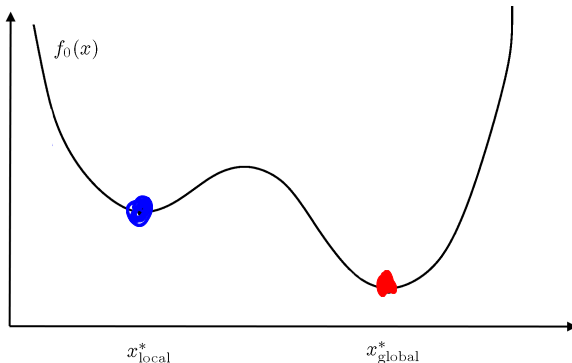
Why?

# Global vs. Local Optimality

## Global optimal point

$\mathbf{x}_{\text{global}}^*$  is a **globally optimal point** if

$$\mathbf{x}_{\text{global}}^* = \arg \min_{\mathbf{x} \in \mathbb{R}^d} f_0(\mathbf{x}) \quad \text{s.t.} \quad \begin{aligned} f_i(\mathbf{x}) &\leq 0, \quad i = 1, \dots, n \\ g_j(\mathbf{x}) &= 0, \quad j = 1, \dots, m \end{aligned}$$

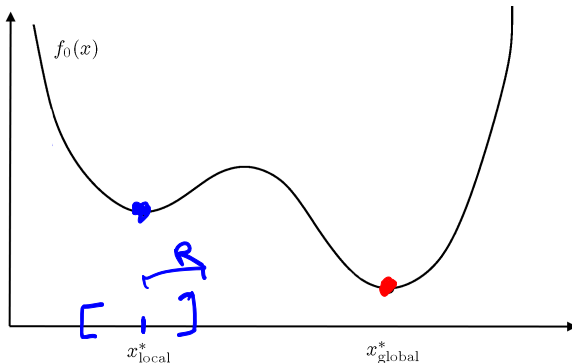


# Global vs. Local Optimality

## Locally optimal point

$\mathbf{x}_{\text{local}}^*$  is a **locally optimal point** if for some  $R > 0$ :

$$\mathbf{x}_{\text{local}}^* = \arg \min_{\mathbf{x} \in \mathbb{R}^d} f_0(\mathbf{x}) \quad \text{s.t.} \quad \begin{aligned} f_i(\mathbf{x}) &\leq 0, \quad i = 1, \dots, n \\ g_j(\mathbf{x}) &= 0, \quad j = 1, \dots, m \\ \|\mathbf{x} - \mathbf{x}_{\text{local}}^*\| &\leq R \end{aligned}$$

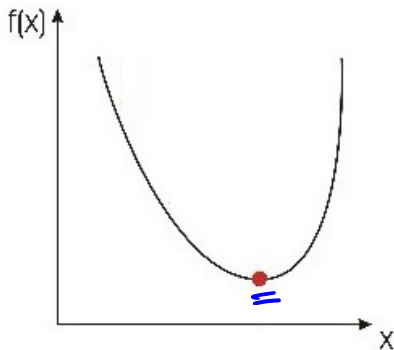


# Why are Convex OPs **easy** to Solve?

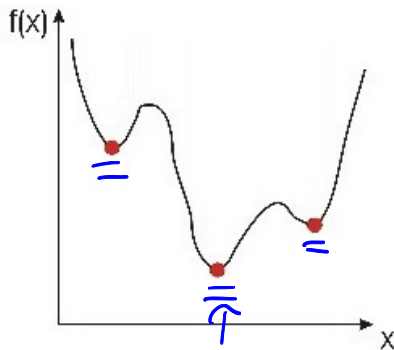
# Why are Convex OPs **easy** to Solve?

## Theorem

Every **locally** optimal point of a convex optimization problem is also **globally** optimal.



convex



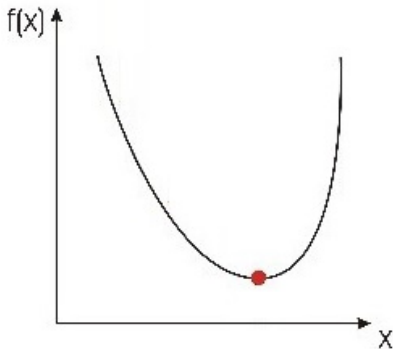
non-convex



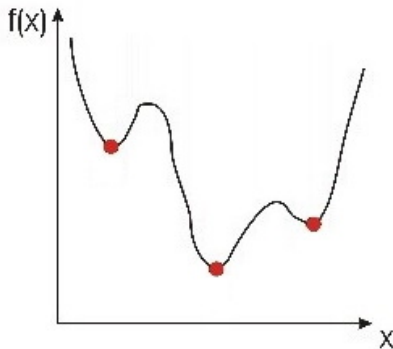
# Why are Convex OPs **easy** to Solve?

## Theorem

Every **locally** optimal point of a convex optimization problem is also **globally** optimal.



convex



non-convex

How can we exploit this property for solving the OP?

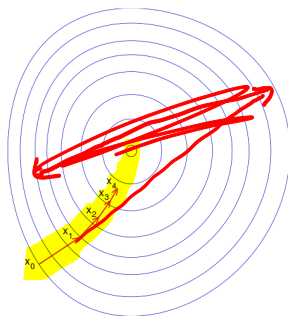
# Gradient Descent

For unconstrained problems, we can iteratively move into **direction of steepest descent** (negative gradient direction) of the objective function:

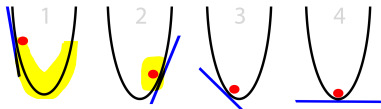
## Gradient descent algorithm

```
1: initialize  $\mathbf{x}_0$  (e.g., randomly)
2: for  $t = 1 : T$  do
3:    $\mathbf{x}_{t+1} := \mathbf{x}_t - \lambda_t \nabla f_0(\mathbf{x}_t)$ 
4: end for
```

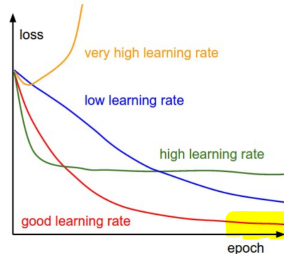
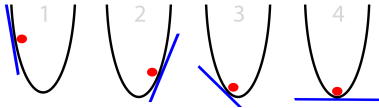
$\lambda_t$  is called **step size** or **learning rate**



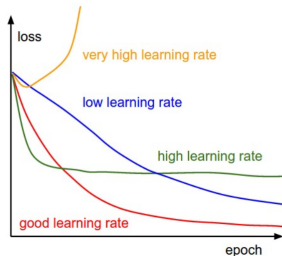
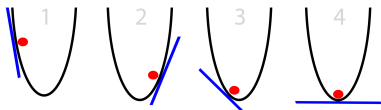
# An Adequate Learning Rate can be Crucial...



# An Adequate Learning Rate can be Crucial...



# An Adequate Learning Rate can be Crucial...



► A typical choice:  $\lambda_t := \frac{1}{t}$

# Gradient Descent Convergences

## Theorem

(Bertsekas, Prop. 1.2.4 & 1.3.3)

Let  $f_0 : \mathbb{R}^d \rightarrow \mathbb{R}$  be an arbitrary (possibly non-convex) objective function. Then, under some assumptions,<sup>1</sup> we have:

- 1 Gradient descent converges.<sup>2</sup>
- 2 For ideal choice of the learning rate, the convergence rate is at least as good as:

$$f_0(\mathbf{x}_t) - f_0(\mathbf{x}_{\text{local}}^*) \leq O(1/t).$$

---

<sup>1</sup> The theorem assumes Lipschitz-continuous gradients with a uniformly bounded Lipschitz constant:  $\exists L : \|\nabla f(\mathbf{x}) - \nabla f(\tilde{\mathbf{x}})\| \leq L \|\mathbf{x} - \tilde{\mathbf{x}}\|$ . Convergence is guaranteed for all learning rate schedules satisfying  $\sum_{t=1}^{\infty} \lambda_t = \infty$  and  $\lambda_t \xrightarrow[t \rightarrow \infty]{} 0$ , but with varying rates of convergence. The favorable  $O(1/t)$  rate is achieved using the minimization rule:  $\lambda_t := \arg \min_{\lambda} f_0(\mathbf{x}_t - \lambda \nabla f_0(\mathbf{x}_t))$ .

<sup>2</sup> More precisely, it converges to a stationary point (that is, a point where the gradient is zero, i.e., either a minimum, a maximum, or a saddle point). However, machine-learning practice (e.g., in deep learning) has shown that this is usually a minimum, so we are good. :)

## Let's Try Applying **Gradient Descent** to the SVM OP:

$$\begin{aligned} \min_{b \in \mathbb{R}, \mathbf{w} \in \mathbb{R}^d, \xi \in \mathbb{R}^n} \quad & \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & \forall i : 1 - \xi_i - y_i(\mathbf{w}^\top \mathbf{x}_i + b) \leq 0, \quad -\xi_i \leq 0 \end{aligned} \quad (\text{SVM})$$

## Let's Try Applying **Gradient Descent** to the SVM OP:

$$\min_{b \in \mathbb{R}, \mathbf{w} \in \mathbb{R}^d, \xi \in \mathbb{R}^n} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i \quad (\text{SVM})$$

$$\text{s.t. } \forall i : 1 - \xi_i - y_i(\mathbf{w}^\top \mathbf{x}_i + b) \leq 0, \quad -\xi_i \leq 0$$

What could be a problem?



## Let's Try Applying **Gradient Descent** to the SVM OP:

$$\min_{b \in \mathbb{R}, \mathbf{w} \in \mathbb{R}^d, \xi \in \mathbb{R}^n} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i \quad (\text{SVM})$$

$$\text{s.t. } \forall i : 1 - \xi_i - y_i(\mathbf{w}^\top \mathbf{x}_i + b) \leq 0, \quad -\xi_i \leq 0$$

What could be a problem? The constraints!

# Let's Try Applying **Gradient Descent** to the SVM OP:

$$\min_{b \in \mathbb{R}, \mathbf{w} \in \mathbb{R}^d, \xi \in \mathbb{R}^n} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \max(0, 1 - \gamma_i (\mathbf{w}^\top \mathbf{x}_i + b)) \quad (\text{SVM})$$

$$\text{s.t. } \forall i: 1 - \xi_i - y_i (\mathbf{w}^\top \mathbf{x}_i + b) \leq 0, \quad -\xi_i \leq 0$$

What could be a problem? The constraints!

## Proposition

The linear SVM can be equivalently re-written as follows:

## Unconstrained linear soft-margin SVM

$$\min_{b \in \mathbb{R}, \mathbf{w} \in \mathbb{R}^d} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \max(0, 1 - y_i (\mathbf{w}^\top \mathbf{x}_i + b))$$

$$\xi_i \geq 1 - \gamma_i (\mathbf{w}^\top \mathbf{x}_i + b) \iff \xi_i \geq \max(0, 1 - \gamma_i (\mathbf{w}^\top \mathbf{x}_i + b))$$

$$\xi_i \geq 0$$

## Let's Try Again Applying Gradient Descent to SVM

The new, unconstrained objective is :

$$f_0(\mathbf{w}, b) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \max(0, 1 - y_i(\mathbf{w}^\top \mathbf{x}_i + b)).$$

## Let's Try Again Applying Gradient Descent to SVM

The new, unconstrained objective is :

$$f_0(\mathbf{w}, b) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \max \left( 0, 1 - y_i(\mathbf{w}^\top \mathbf{x}_i + b) \right).$$

What could be another problem?

# Let's Try Again Applying Gradient Descent to SVM

The new, unconstrained objective is :

$$f_0(\mathbf{w}, b) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \max \left( 0, 1 - \overbrace{y_i(\mathbf{w}^\top \mathbf{x}_i + b)}^{=: x} \right).$$

What could be another problem?  $f_0$  is not differentiable!

$$f: \mathbb{R} \rightarrow \mathbb{R} \\ x \mapsto \max(0, 1 - x)$$



## Let's Try Again Applying Gradient Descent to SVM

The new, unconstrained objective is :

$$f_0(\mathbf{w}, b) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \max \left( 0, 1 - y_i(\mathbf{w}^\top \mathbf{x}_i + b) \right).$$

What could be another problem?  $f_0$  is not differentiable!

We have two options to address this problem:

# Let's Try Again Applying Gradient Descent to SVM

The new, unconstrained objective is :

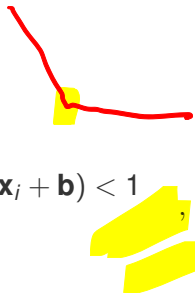
$$f_0(\mathbf{w}, b) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \max(0, 1 - y_i(\mathbf{w}^\top \mathbf{x}_i + b)).$$

What could be another problem?  $f_0$  is not differentiable!

We have two options to address this problem:

- Option 1: We can consider the **subgradient**,

$$\begin{aligned} \nabla \max(0, 1 - y_i(\mathbf{w}^\top \mathbf{x}_i + b)) \\ := \begin{cases} \nabla(1 - y_i(\mathbf{w}^\top \mathbf{x}_i + b)) & \text{if } y_i(\mathbf{w}^\top \mathbf{x}_i + b) < 1 \\ 0 & \text{elsewise} \end{cases} \end{aligned}$$



# Let's Try Again Applying **Gradient Descent** to SVM

The new, unconstrained objective is :

$$f_0(\mathbf{w}, b) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \max(0, 1 - y_i(\mathbf{w}^\top \mathbf{x}_i + b)).$$

What could be another problem?  $f_0$  is not differentiable!

We have two options to address this problem:

- Option 1: We can consider the **subgradient**,

$$\begin{aligned} \sum_{i=1}^n \nabla \max(0, 1 - y_i(\mathbf{w}^\top \mathbf{x}_i + b)) \\ := \begin{cases} \sum_{i=1}^n \nabla(1 - y_i(\mathbf{w}^\top \mathbf{x}_i + b)) & \text{if } y_i(\mathbf{w}^\top \mathbf{x}_i + b) < 1 \\ 0 & \text{elsewise} \end{cases}, \end{aligned}$$

and then use **subgradient descent**: that is, gradient descent but using the subgradient in place of the gradient.



# Logistic Regression

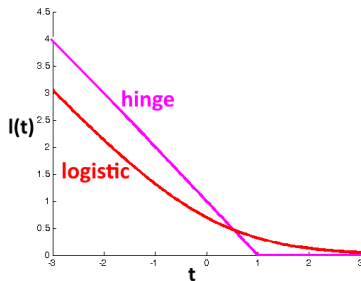
Option 2:

- ▶ Replace 'hinge' function

$$f(x) = \max(0, 1 - x)$$

appearing in SVM by its smooth approximation, the 'logistic' function:

$$l(x) = -\ln(1 + \exp(-x)).$$



This results in:

## Logistic Regression (LR)

$$\min_{b \in \mathbb{R}, \mathbf{w} \in \mathbb{R}^d} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \ln(1 + \exp(-y_i(\mathbf{w}^\top \mathbf{x}_i + b)))$$

LR is differentiable, so we can solve it using gradient descent.

# But (Sub)gradient Descent of SVM / LR is Too Slow!

# But (Sub)gradient Descent of SVM / LR is Too Slow!

- ▶ Each evaluation of the subgradient involves a FOR loop over all data points:  
FOR all  $i = 1, \dots, n$
- ▶ This is costly ( $n$  large for big data)!

# But (Sub)gradient Descent of SVM / LR is Too Slow!

- ▶ Each evaluation of the subgradient involves a FOR loop over all data points:  
FOR all  $i = 1, \dots, n$
- ▶ This is costly ( $n$  large for big data)!

But evaluating the full FOR loop can be unnecessary:

# But (Sub)gradient Descent of SVM / LR is Too Slow!

- ▶ Each evaluation of the subgradient involves a FOR loop over all data points:  
FOR all  $i = 1, \dots, n$
- ▶ This is costly ( $n$  large for big data)!

But evaluating the full FOR loop can be unnecessary:

- ▶ Imagine all data points were the same (duplicates), thus:

$n \times$

# But (Sub)gradient Descent of SVM / LR is Too Slow!

- ▶ Each evaluation of the subgradient involves a FOR loop over all data points:  
FOR all  $i = 1, \dots, n$
- ▶ This is costly ( $n$  large for big data)!

But evaluating the full FOR loop can be unnecessary:

- ▶ Imagine all data points were the same (duplicates), thus:
- ▶  $\sum_{i=1}^n \max(0, 1 - y_i(\mathbf{w}^\top \mathbf{x}_i + b)) = n \cdot \max(0, 1 - y_1(\mathbf{w}^\top \mathbf{x}_1 + b))$

# But (Sub)gradient Descent of SVM / LR is Too Slow!

- ▶ Each evaluation of the subgradient involves a FOR loop over all data points:  
FOR all  $i = 1, \dots, n$
- ▶ This is costly ( $n$  large for big data)!

But evaluating the full FOR loop can be unnecessary:

- ▶ Imagine all data points were the same (duplicates), thus:
- ▶  $\sum_{i=1}^n \max(0, 1 - y_i(\mathbf{w}^\top \mathbf{x}_i + b)) = n \cdot \max(0, 1 - y_1(\mathbf{w}^\top \mathbf{x}_1 + b))$

Real data is usually less extreme (not exact duplicates),  
but yet contains lots of redundancy

⇒ unnecessary to evaluate the full subgradient in every  
iteration



# Stochastic Gradient Descent (SGD) is Much Faster...

## Stochastic subgradient descent algorithm (SVM)

```
1: initialize  $(b, \mathbf{w})_0$  (e.g., randomly)
2: for  $t = 1 : T$  do
3:   Randomly select  $B$  many data points
4:   Denote their indexes by  $I \subset \{1, \dots, n\}$  (i.e.,  $|I| = B$ )
5:   Update  $(b, \mathbf{w})_{t+1} :=$   

   
$$(b, \mathbf{w})_t - \lambda_t \nabla \left( \frac{1}{2} \|\mathbf{w}\|^2 + \frac{Cn}{B} \sum_{i \in I} \max(0, 1 - y_i(\mathbf{w}^\top \mathbf{x}_i + b)) \right)$$

6: end for
```


The classic SGD algorithm uses just a single data point per iteration ( $B = 1$ ). Nowadays more common in ML: **mini-batch SGD**, where we use an intermediate value, such as  $B = 100$ .



# Stochastic Gradient Descent (SGD) is Much Faster...

## Stochastic subgradient descent algorithm (SVM)

- 1: initialize  $(b, \mathbf{w})_0$  (e.g., randomly)
- 2: **for**  $t = 1 : T$  **do**
- 3:     Randomly select  $B$  many data points
- 4:     Denote their indexes by  $I \subset \{1, \dots, n\}$  (i.e.,  $|I| = B$ )
- 5:     Update  $(b, \mathbf{w})_{t+1} :=$

$$(b, \mathbf{w})_t - \lambda_t \nabla \left( \frac{1}{2} \|\mathbf{w}\|^2 + \frac{Cn}{B} \sum_{i \in I} \max(0, 1 - y_i(\mathbf{w}^\top \mathbf{x}_i + b)) \right)$$


- 6: **end for**

The **batch size**  $B \in [1, n]$  needs to be chosen a priori.

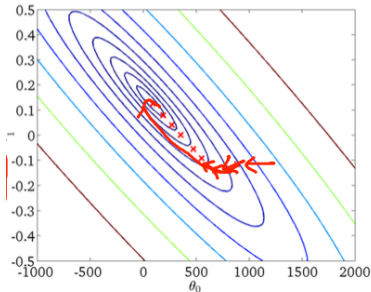
For LR simply replace in Line 5 the term

$\max(0, 1 - y_i(\mathbf{w}^\top \mathbf{x}_i + b))$  by  $\log(1 + \exp(-y_i(\mathbf{w}^\top \mathbf{x}_i + b)))$ .

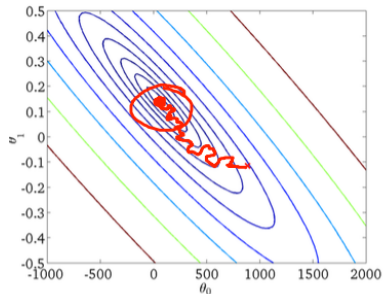
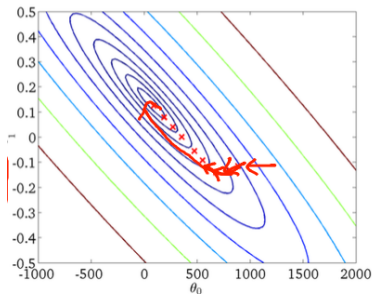
---

The classic SGD algorithm uses just a single data point per iteration ( $B = 1$ ). Nowadays more common in ML: **mini-batch SGD**, where we use an intermediate value, such as  $B = 100$ .

# What is the Difference Between Gradient Descent and Stochastic Gradient Descent?



# What is the Difference Between Gradient Descent and Stochastic Gradient Descent?



# SGD Convergences

## Theorem

(e.g., Bottou et al., 2018)

Consider SGD using the learning rate  $\lambda_t := 1/t$ . Then, under mild assumptions, SGD converges with high probability to a stationary point<sup>1</sup> with rate:

$$f_0(\mathbf{x}_t) - f_0(\mathbf{x}_{\text{local}}^*) \leq O(1/t)$$

Remark: this holds also for non-convex  $f_0$

---

<sup>1</sup> In ML practice, this is usually a local minimum.

# Software

An extremely fast implementation of SGD applied to the linear (soft-margin) SVM and logistic regression is contained in **Vowpal Wabbit**<sup>1</sup>

- ▶ Industry standard in ad prediction

Discussion:

- ▶ VW trades speed for accuracy, so use it only when really necessary (=big data)
- ▶ Otherwise use LIBLINEAR, a more accurate and still relatively fast solver for SVM and LR<sup>2</sup>

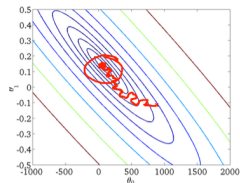
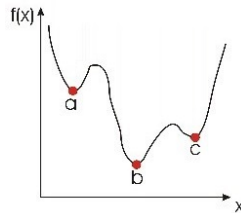
<sup>1</sup> [https://github.com/VowpalWabbit/vowpal\\_wabbit/wiki](https://github.com/VowpalWabbit/vowpal_wabbit/wiki)

<sup>2</sup> <https://www.csie.ntu.edu.tw/~cjlin/liblinear/>



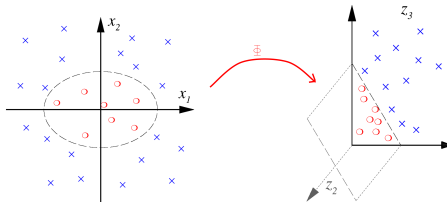
# Conclusion

- ▶ Convex optimization problem (OP):
  - ▶ optimize convex function over convex set
- ▶ Why convexity?
  - ▶ cannot get stuck in local optimum
- ▶ SVM can be formulated as convex OP
  - ▶ unconstrained formulation of SVM
  - ▶ gradient descent slow for SVM
  - ▶ solution: stochastic gradient descent (Vowpal Wabbit)



# Next Week

- ▶ Non-linear SVM
- ▶ Kernel methods



# Refs I



S. Shalev-Shwartz, Y. Singer, N. Srebro, and A. Cotter, Pegasos: Primal estimated sub-gradient solver for svm, *Mathematical programming*, vol. 127, no. 1, pp. 3–30, 2011.



L. Bottou, F. E. Curtis, and J. Nocedal, Optimization methods for large-scale machine learning, *Siam Review*, vol. 60, no. 2, pp. 223–311, 2018.