

FrankWolfe.jl: A high-performance and flexible toolbox for Frank-Wolfe algorithms and Conditional Gradients

Mathieu Besançon

Zuse Institute Berlin, Germany

BESANCON@ZIB.DE

Alejandro Carderera

*Georgia Institute of Technology, USA
Zuse Institute Berlin, Germany*

ALEJANDRO.CARDERERA@GATECH.EDU

Sebastian Pokutta

*Technische Universität Berlin, Germany
Zuse Institute Berlin, Germany*

POKUTTA@MATH.TU-BERLIN.DE

Abstract

We present `FrankWolfe.jl`, an open-source implementation of several popular Frank-Wolfe and Conditional Gradients variants for first-order constrained optimization. The package is designed with flexibility and high-performance in mind, allowing for easy extension and relying on few assumptions regarding the user-provided functions. It supports Julia's unique multiple dispatch feature, and interfaces smoothly with generic linear optimization formulations using `MathOptInterface.jl`.

1. Introduction

We provide an open-source software package released under the MIT license, for *Frank-Wolfe (FW)* and *Conditional Gradients* algorithms implemented in Julia (Bezanson et al., 2017). Its focus is a general class of constrained convex optimization problems of the form:

$$\min_{\mathbf{x} \in C} f(\mathbf{x})$$

where \mathcal{X} is a Hilbert space, $C \subseteq \mathcal{X}$ is a compact convex set, and $f : C \rightarrow \mathbb{R}$ is a convex, continuously differentiable function.

Although the Frank-Wolfe algorithm and its variants have been studied for more than half a decade and have gained a lot of attention for their theoretical and computational properties, no de-facto implementation exists. The goal of the package is to become a reference open-source implementation for practitioners in need of a flexible and efficient first-order method and for researchers developing and comparing new approaches on similar classes of problems.

2. Frank-Wolfe algorithms

The Frank-Wolfe algorithm (Frank & Wolfe, 1956) (also known as the Conditional Gradient algorithm (Levitin & Polyak, 1966)), is a first-order algorithm for constrained optimization that avoids the use of projections at each iteration. For the sake of exposition we confine ourselves to the Euclidean setup. The main ingredients that the FW algorithm leverages are:

1. **First-Order Oracle (FOO)**: Given $\mathbf{x} \in C$, the oracle returns $\nabla f(\mathbf{x})$.
2. **Linear Minimization Oracle (LMO)**: Given $\mathbf{d} \in \mathcal{X}$, the oracle returns:

$$\mathbf{v} \in \operatorname{argmin}_{\mathbf{x} \in C} \langle \mathbf{d}, \mathbf{x} \rangle \tag{1}$$

The simplest version of the algorithm (shown in Algorithm 2.1) builds a linear approximation to the function at a given iterate, using first-order information, and minimizes this approximation over the feasible

region (Line 2 in Algorithm 2.1). A schematic representation of the step is described in Fig. 1 where the blue curves represent the contour lines of $f(\mathbf{x})$, and the red lines represent the contour lines of the linear approximation built at \mathbf{x}_t . The new iterate is then computed as a convex combination of the current iterate and the linear minimizer \mathbf{v}_t from the LMO (Line 3 in Algorithm 2.1). Since \mathbf{v}_t is an extreme point of the feasible region, the new iterate is feasible by convexity of C and remains so throughout the algorithm. Alternatively, one can view each iteration of the Frank-Wolfe algorithm as finding the direction that is best aligned with the negative of the gradient using only the current iterate and the extreme points of the feasible region.

Algorithm 2.1 Frank-Wolfe algorithm

Require: Point $\mathbf{x}_0 \in C$, function f , step-sizes $\gamma_t > 0$.

Ensure: Iterates $\mathbf{x}_1, \dots \in C$.

```

1: for  $t = 0$  to  $\dots$  do
2:    $\mathbf{v}_t \leftarrow \operatorname{argmin}_{\mathbf{v} \in C} \langle \nabla f(\mathbf{x}_t), \mathbf{v} \rangle$ 
3:    $\mathbf{x}_{t+1} \leftarrow \mathbf{x}_t + \gamma_t (\mathbf{v}_t - \mathbf{x}_t)$ 
4: end for
```

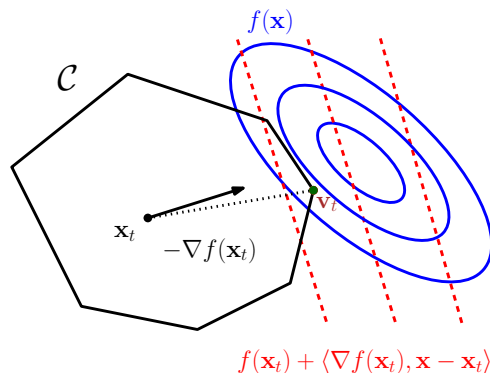


Figure 1: Schematic of a Frank-Wolfe step at iteration t .

The Frank-Wolfe algorithm and its variants have recently gained a lot of attention in the Machine Learning community due to the fact that they avoid the use of projection operators to ensure the feasibility of the iterates. This is particularly advantageous in applications where computing a projection (in essence, solving a quadratic problem) is much more computationally expensive than solving a linear optimization problem (Combettes & Pokutta, 2021). For example, if C is the set of matrices in $\mathbb{R}^{n \times m}$ of bounded nuclear norm, computing the projection of a matrix onto C requires computing a full Singular Value Decomposition. On the other hand, solving a linear optimization problem over C only requires computing the left and right singular vectors associated with the largest singular value.

Another interesting property of the FW algorithm is that when minimizing a convex function f , the primal gap $f(\mathbf{x}) - f(\mathbf{x}^*)$ is bounded by a quantity dubbed the *Frank-Wolfe gap* or the *dual gap* defined as $\langle \nabla f(\mathbf{x}_t), \mathbf{x}_t - \mathbf{v}_t \rangle$. This follows directly from the convexity of f and the maximality of \mathbf{v}_t since:

$$f(\mathbf{x}) - f(\mathbf{x}^*) \leq \langle \nabla f(\mathbf{x}_t), \mathbf{x}_t - \mathbf{x}^* \rangle \leq \max_{\mathbf{v} \in C} \langle \nabla f(\mathbf{x}_t), \mathbf{x}_t - \mathbf{v} \rangle = \langle \nabla f(\mathbf{x}_t), \mathbf{x}_t - \mathbf{v}_t \rangle.$$

The dual gap (and its stronger variants) is extremely useful as a progress and stopping criterion when running Frank-Wolfe algorithms.

2.1 Algorithm variants

The simple (and robust) algorithm presented in Algorithm 2.1 is not able to attain linear convergence in general, which has led to the development of algorithmic variants that enhance the performance of the original algorithm, while maintaining many of its advantageous properties. We summarize below the central ideas of the variants implemented in the package and highlight in Table 1 key properties that can drive the choice of a variant on a given use case. We mention briefly that most variants also work for the nonconvex case, providing some locally optimal solution in this case.

Standard Frank-Wolfe The simplest Frank-Wolfe variant is presented in Algorithm 2.1. It has the lowest memory requirements out of all the variants, as in its simplest form only requires keeping track of the current iterate. As such it is suited for extremely large problems. However, in certain cases, this comes at the cost of speed of convergence in terms of iteration count, when compared to other variants. As an example, when

minimizing a strongly convex and smooth function over a polytope this algorithm might converge sublinearly, whereas the three variants that will be presented next converge linearly.

Away-step Frank-Wolfe The most popular among the Frank-Wolfe variants is the *Away-step Frank-Wolfe* (AFW) algorithm (Guélat & Marcotte, 1986; Lacoste-Julien & Jaggi, 2015). While the FW algorithm presented in Algorithm 2.1 can only move *towards* extreme points of C , the AFW can move *away* from some extreme points of C , hence the name of the algorithm. To be more specific, the AFW algorithm moves away from vertices in its *active set* at iteration t , denoted by \mathcal{S}_t , which contains the set of vertices \mathbf{v}_k for $k < t$ that allow us to recover the current iterate as a convex combination. This algorithm expands the range of directions that the FW algorithm can move along, at the expense of having to explicitly maintain the current iterate as a convex decomposition of extreme points. See Fig. 2 for a schematic of the behavior of the two algorithms for a simple example.

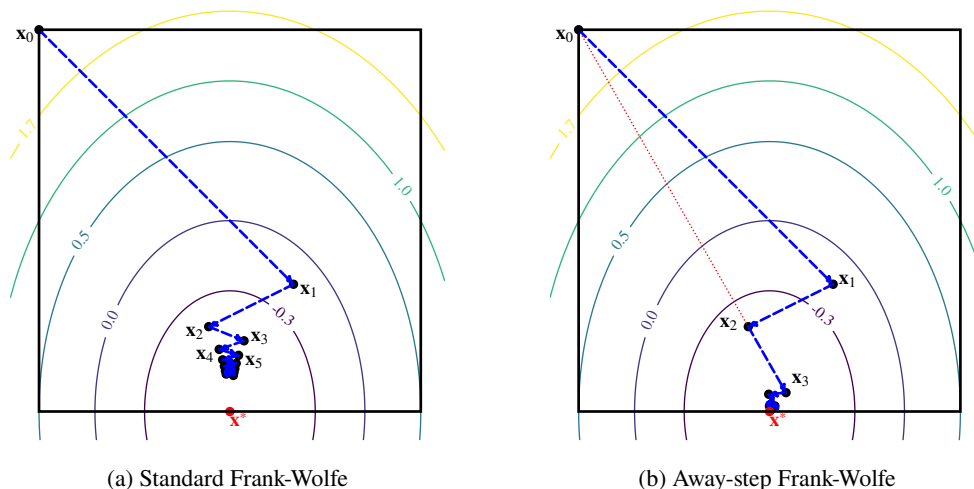


Figure 2: Convergence comparison of the standard Frank-Wolfe algorithm and the Away-step Frank-Wolfe algorithm when minimizing a quadratic function (contour lines depicted) over a simple polytope using exact line search. Note that the minimizer is on a face of the polytope, which causes the Frank-Wolfe algorithm to slowly *zig-zag* towards the solution, causing sublinear convergence

Lazifying Frank-Wolfe variants One running assumption for the two previous variants is that calling the LMO is cheap. There are many applications where calling the LMO in absolute terms is costly (but is cheap in relative terms when compared to performing a projection). In such cases, one can attempt to *lazify* FW algorithms, to avoid having to compute $\operatorname{argmin}_{\mathbf{v} \in C} \langle \nabla f(\mathbf{x}_t), \mathbf{v} \rangle$ by calling the LMO, settling for finding approximate minimizers that guarantee enough progress (Braun et al., 2017). This allows us to substitute the LMO by a *Weak Separation Oracle* while maintaining essentially the same convergence rates. In practice, these algorithms search for appropriate vertices among the vertices in a cache, or the vertices in the active set \mathcal{S}_t , and can be much faster in wall-clock time. In the package, both AFW and FW have lazy variants while the BCG algorithm is lazified by design.

Blended Conditional Gradients The FW and AFW algorithms, and their lazy variants share one feature: they attempt to make primal progress over a reduced set of vertices. The AFW algorithm does this through away steps (which do not increase the cardinality of the active set), and the lazy variants do this through the use of previously exploited vertices. A third strategy that one can follow is to explicitly *blend* Frank-Wolfe steps with gradient descent steps over the convex hull of the active set (note that this can be done without requiring a projection oracle over C , thus making the algorithm projection-free). This results in the *Blended*

Conditional Gradient (BCG) algorithm (Braun et al., 2019), which attempts to make as much progress as possible through the convex hull of the current active set S_i until it automatically detects that in order to further make further progress it requires additional calls to the LMO.

Stochastic Frank-Wolfe In many problem instances, evaluating the FOO at a given point is prohibitively expensive. In such cases, one usually has access to a *Stochastic First-Order Oracle* (SFOO), from which one can build a gradient estimator (Robbins & Monro, 1951). This idea, which has powered much of the success of deep learning, can also be applied to the Frank-Wolfe algorithm (Hazan & Luo, 2016), resulting in the *Stochastic Frank-Wolfe* (SFW) algorithm and its variants.

Algorithm	Convergence		Sparsity	Numerical Stability	Active Set?	Lazifiable?
	Progress/iteration	Time/iteration				
FW	Low	Low	Low	High	No	Yes
AFW	Medium	Medium-High	Medium	Medium-High	Yes	Yes
BCG	High	Medium-High	High	Medium	Yes	By design
SFW	Low	Low	Low	High	No	No

Table 1: Comparison of the characteristics of the algorithms in the package, when applied to a general problem.

3. The `FrankWolfe.jl` package

The package offers the first implementation of Frank-Wolfe variants in Julia but more broadly tackles constrained optimization problems in a way that complements well the thriving Julia optimization ecosystem.

Unlike disciplined convex frameworks or algebraic modeling languages such as `Convex.jl` (Udell et al., 2014) or `JuMP.jl` (Dunning et al., 2017; Legat et al., 2020), our framework allows for arbitrary Julia functions defined outside of a Domain-Specific Language. Users can provide their gradient implementation or leverage one of the many automatic differentiation packages available in Julia.

Several ecosystems have emerged for non-linear optimization, but these often only handle specific types of constraints. Most notably, `JSOSolvers.jl` (Orban et al., 2020) offers a trust-region second-order method for bound constraints. On the other hand, `Optim.jl` (Mogensen & Riseth, 2018) implements an interior-point method handling bound and non-linear constraints. Constrained first-order methods based on proximal operators are implemented in `StructuredOptimization.jl` and `ProximalAlgorithms.jl` (Antonello et al., 2018) but only allow specific functions defined through the exposed Domain-Specific Language.

One central design principle of `FrankWolfe.jl` is to rely on few assumptions regarding the user-provided functions, the atoms returned by the LMO, and their implementation. The package works for instance out of the box when the LMO returns Julia subtypes of `AbstractArray`, representing finite-dimensional vectors, matrices or higher-order arrays.

Another design principle has been to favor in-place operations and reduce memory allocations when possible, since these can become expensive when repeated at all iterations. This is reflected in the memory emphasis mode (the default mode for all algorithms), where as many computations as possible are performed in-place as well as in the gradient interface, where the gradient function is provided with a variable to write into rather than reallocating every time a gradient is computed. The performance difference can be quite pronounced for problems in large dimensions, for example passing a gradient of size 7.5GB on a state-of-the-art machine is about 8 times slower than an in-place update.

Finally, default parameters are chosen to make all algorithms as robust as possible out of the box, while allowing extension and fine tuning for advanced users. For example, the default step size strategy for all (but the stochastic variant) is the adaptive step size rule of Pedregosa et al. (2020), which in computations not only usually outperforms both line search and the short step rule by dynamically estimating the Lipschitz

constant but also overcomes several issues with the limited additive accuracy of traditional line search rules. Similarly, the BCG variant automatically upgrades the numerical precision for certain subroutines if numerical instabilities are detected.

3.1 Linear minimization oracle interface

One key step of FW algorithms is the linear minimization step which, given first-order information at the current iterate, returns a vertex of the feasible region that minimizes the linear approximation of the function. It is defined in `FrankWolfe.jl` using a single function:

```
function compute_extreme_point(lmo::LMO, direction::D; kwargs...):V
    # ...
end
```

The first argument `lmo` represents the linear minimization oracle for the specific problem. It encodes the feasible region C , but also some algorithmic parameters or state. This is especially useful for the lazified FW variants, as in these cases the LMO types can take advantage of caching, by storing the extreme vertices that have been computed in previous iterations and then looking up vertices from the cache before computing a new one.

The package implements LMOs for commonly encountered feasible regions including L_p -norm balls, K -sparse polytopes, the Birkhoff polytope, and the nuclear norm ball for matrix spaces, leveraging known closed forms of extreme points. The multiple dispatch mechanism allows for different implementations of a single LMO with multiple direction types. The type `V` used to represent the computed vertex is also specialized to leverage the properties of extreme vertices of the feasible region. For instance, although the Birkhoff polytope is the convex hull of all doubly stochastic matrices of a given dimension, its extreme vertices are permutation matrices that are much sparser in nature. We also leverage sparsity outside of the traditional sense of nonzero entries. When the feasible region is the nuclear norm ball in $\mathbb{R}^{N \times M}$, the vertices are rank-one matrices. Even though these vertices are dense, they can be represented as the outer product of two vectors and thus be stored with $\mathcal{O}(N + M)$ entries instead of $\mathcal{O}(N \times M)$ for the equivalent dense matrix representation. The Julia abstract matrix representation allows the user and the library to interact with these rank-one matrices with the same API as standard dense and sparse matrices.

In some cases, users may want to define a custom feasible region that does not admit a closed-form linear minimization solution. We implement a generic LMO based on `MathOptInterface.jl` (Legat et al., 2020), thus allowing users on the one hand to select any off-the-shelf LP, MILP, or conic solver suitable for their problem, and on the other hand to formulate the constraints of the feasible domain using the `JuMP.jl` or `Convex.jl` DSL. Furthermore, the interface is naturally extensible by users who can define their own LMO and implement the corresponding `compute_extreme_point` method.

3.2 Numeric type genericity

The package was designed from the start to be generic over both the used numeric types and data structures. Numeric type genericity allows running the algorithms in extended fixed or arbitrary precision, e.g., the package works out-of-the-box with `Double64` and `BigFloat` types. Extended precision is essential for high-dimensional problems where the condition number of the gradient computation becomes too high. For some well-conditioned problems, reduced precision is sometimes sufficient to achieve the desired tolerance. Furthermore, it opens the possibility of gradient computation and LMO steps on hardware accelerators such as GPUs.

4. Application examples

We will now present a few examples that highlight specific features of the package. The full code of each example (and several more) can be found in the examples folder of the repository.

4.1 Polynomial Regression

The data encountered in many real-world applications can be modeled as coming from data-generating processes, which non-linearly map a set of input features to an output space. In many cases, this non-linear mapping can be modeled with a relatively sparse linear combination of simpler non-linear functions, which we call basis functions. For example, if we denote the scalar output by $y \in \mathbb{R}$, the input features by $\mathbf{x} \in \mathbb{R}^n$, and the library of basis functions by $f_i: \mathbb{R}^n \rightarrow \mathbb{R}$ with $i \in [1, m]$ we might have that the data is generated as:

$$y = \sum_{i=1}^m c_i f_i(\mathbf{x}),$$

where $c_i \in \mathbb{R}$, and only a few of the basis functions participate in the data generating process, i.e., many of the coefficients satisfy $c_i = 0$. Our task is then to recover which of these c_i coefficients are non-zero, along with their corresponding value, when we are given a collection of input/output measurements $\{\mathbf{x}_j, y_j\}_{j=1}^N$. In the absence of noise this problem can be written as:

$$\begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix} = \begin{bmatrix} f_1(\mathbf{x}_1) & \cdots & f_m(\mathbf{x}_1) \\ \vdots & \ddots & \vdots \\ f_1(\mathbf{x}_N) & \cdots & f_m(\mathbf{x}_N) \end{bmatrix} \begin{bmatrix} c_1 \\ \vdots \\ c_m \end{bmatrix}, \quad (2)$$

which we can succinctly write as $\mathbf{y} = \mathbf{A}\mathbf{c}$ by collecting the points $\{\mathbf{x}_j\}_{j=1}^N$ into a vector $\mathbf{y} \in \mathbb{R}^N$, the coefficients $\{c_j\}_{j=1}^m$ into a vector $\mathbf{c} \in \mathbb{R}^m$, and the basis function evaluations into an appropriate matrix \mathbf{A} . If the rank of the matrix \mathbf{A} is equal to the number of unknowns (in this case m), a direct linear solve provides the unique optimal solution. On the other hand, in the case where $N < m$, the uniqueness of the optimal solution does not hold. In order for the previous case to make sense, we need to impose further structure on the problem in Equation (2). This is where the sparsity assumption comes into play, as we know that only a few of the c_i coefficients are non-zero.

If we were to select an appropriate subset of the basis functions such that the matrix of this *reduced* problem has a rank equal to the *reduced* number of unknowns, we could again solve exactly the linear system to find these non-zero coefficients. The problem would then become that of finding the *reduced* subset of nonzero coefficients that perfectly reproduces the data generation. That is, we would be trying to solve:

$$\min_{\mathbf{y}=\mathbf{A}\mathbf{c}} \|\mathbf{c}\|_0, \quad (3)$$

where $\|\mathbf{c}\|_0$ counts the non-zero entries in \mathbf{c} . Unfortunately, this combinatorial problem is NP-hard (see [Juditsky & Nemirovski \(2020\)](#) for an excellent discussion). To add insult to injury, we typically deal with data that is contaminated with noise, thus turning the problem into that of finding an estimator of the parameters. We can *convexify* the problem by tackling a relaxation of the form:

$$\min_{\mathbf{c} \in C} \|\mathbf{y} - \mathbf{A}\mathbf{c}\|^2, \quad (4)$$

where C is a convex set that can be used to control the structure of the coefficients; a typical choice is the ℓ_1 norm ball. One of the benefits of solving Problem (4), apart from the fact that it can be approached with the tools of convex optimization, is that there exists a rich literature on the guarantees that we can expect from this recovery (see e.g., [Donoho \(2006\)](#); [Candès & Wakin \(2008\)](#); [Candès et al. \(2006\)](#); [Tropp \(2006\)](#)).

We will present a finite-dimensional example where the basis functions are monomials of the features of the vector $\mathbf{x} \in \mathbb{R}^{15}$ of maximum degree d , that is $f_i(\mathbf{x}) = \prod_{j=1}^n x_j^{a_j}$ with $a_j \in \mathbb{N}$ and $\sum_{j=1}^n a_j \leq d$. We generate a random vector \mathbf{c} that will have 5% non-zero entries drawn from a normal distribution with mean equal to 10 and unit variance. In order to evaluate the polynomial, we generate a total of 1000 data points $\{\mathbf{x}_i\}_{i=1}^N$ from the standard multivariate Gaussian in \mathbb{R}^{15} , with which we will compute the output variables $\{y_i\}_{i=1}^N$. Before evaluating the polynomial these points will be contaminated with noise drawn from a standard multivariate Gaussian. For such a low number of features in the extended space, the polynomial features

can be precomputed, thus reducing the computational burden of each function and gradient evaluation. We leverage `MultivariatePolynomials.jl` (Legat et al., 2021) to create the input polynomial of degree up to 4 in \mathbb{R}^{15} and evaluate it on the training and test data.

The sparsity inducing convex feasible we minimize over is the ℓ_1 norm ball. Solving a linear minimization problem over this feasible region generates points with only one non-zero element. Moreover, there is a closed-form solution for these minimizers. We run the Lazy Away-Step Frank Wolfe (L-AFW) and BCG algorithms with the adaptive line search strategy from Pedregosa et al. (2020), and compare them to Projected Gradient Descent using a smoothness estimate. We will evaluate the output solution on test points drawn in a similar manner as the training points. The radius of the ℓ_1 norm ball that we will use to regularize the problem will be equal to $0.95 \|\mathbf{c}\|_1$.

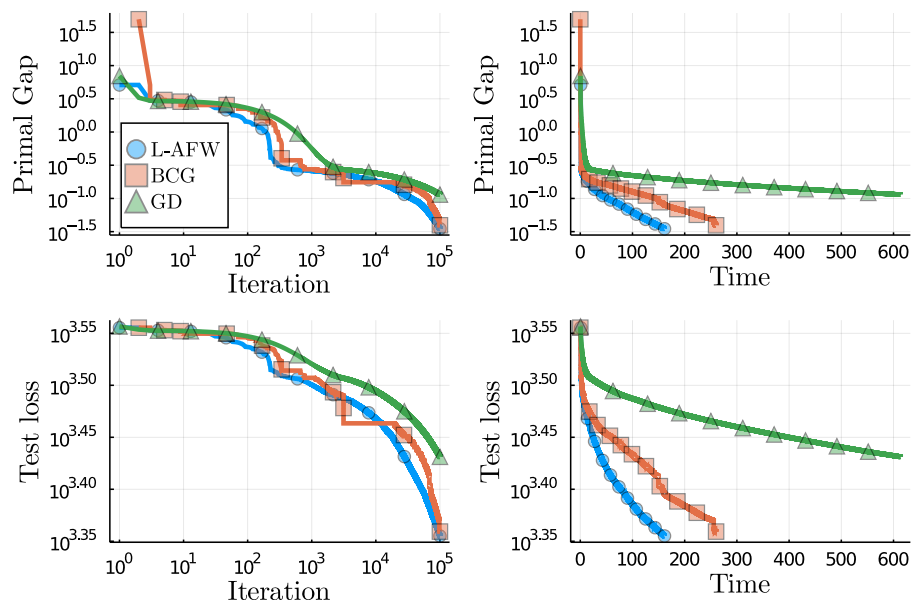


Figure 3: Polynomial regression results.

The primal gap and the test error evolution during the run of the algorithms are presented in Section 4.1, in terms of iteration count and wall-clock time. The L-AFW and BCG algorithms decrease the primal gap faster than the GD algorithm in terms of iteration count, due to the effectiveness of the adaptive stepsize strategy. This faster convergence in iteration count is amplified with respect to time, making the L-AFW and BCG algorithms converge much faster than GD, due to the fact that the L-AFW and BCG algorithms exploit the inherent sparsity of the vertices of the ℓ_1 ball, making the computational cost of each iteration much lower than the cost of each GD step. Note that none of the algorithms manage to perfectly recover the support of the exact coefficients of \mathbf{c} , due to the magnitude of the noise introduced in the example. This is expected, given the small values that some of these coefficients have.

4.2 Matrix completion

Missing data imputation is a key topic in data science. Given a set of observed entries from a matrix $Y \in \mathbb{R}^{m \times n}$, we want to compute a matrix $X \in \mathbb{R}^{m \times n}$ that minimizes the sum of squared errors on the observed entries. As it stands this problem formulation is not well-defined or useful, as one could minimize the objective function simply by setting the observed entries of X to match those of Y , and setting the remaining entries of X arbitrarily. However, this would not result in any meaningful information regarding the unobserved entries in Y , which is one of the key tasks in missing data imputation. A common way to solve this problem is to reduce the degrees of freedom of the problem in order to recover the matrix Y from a small subset of its entries, e.g.,

by assuming that the matrix Y has low rank (Candès & Recht, 2009; Candès & Tao, 2010; Candès & Plan, 2010; Udell & Townsend, 2019). Note that even though the matrix Y has $m \times n$ coefficients, if it has rank r , it can be expressed using only $(m + n - r)r$ coefficients through its singular value decomposition. Finding the matrix $X \in \mathbb{R}^{m \times n}$ with minimum rank whose observed entries are equal to those of Y is a non-convex problem that is $\exists\mathbb{R}$ -hard (Bertsimas et al., 2020). A common proxy for rank constraints is the use of constraints on the nuclear norm of a matrix, which is equal to the sum of its singular values, and can model the convex envelope of matrices of a given rank (Fazel, 2002). Using this property, one of the most common ways to tackle matrix completion problems is to solve:

$$\min_{\|X\|_* \leq \tau} \sum_{(i,j) \in \mathcal{I}} (X_{i,j} - Y_{i,j})^2, \quad (5)$$

where $\tau > 0$ and \mathcal{I} denotes the indices of the observed entries of Y . In this example, we compare the Frank-Wolfe implementation from the package with a Projected Gradient Descent (PGD) algorithm which, after each gradient descent step, projects the iterates back onto the nuclear norm ball. We use one of the movielens datasets to compare the two methods. The code required to reproduce the full example is presented in Appendix A.

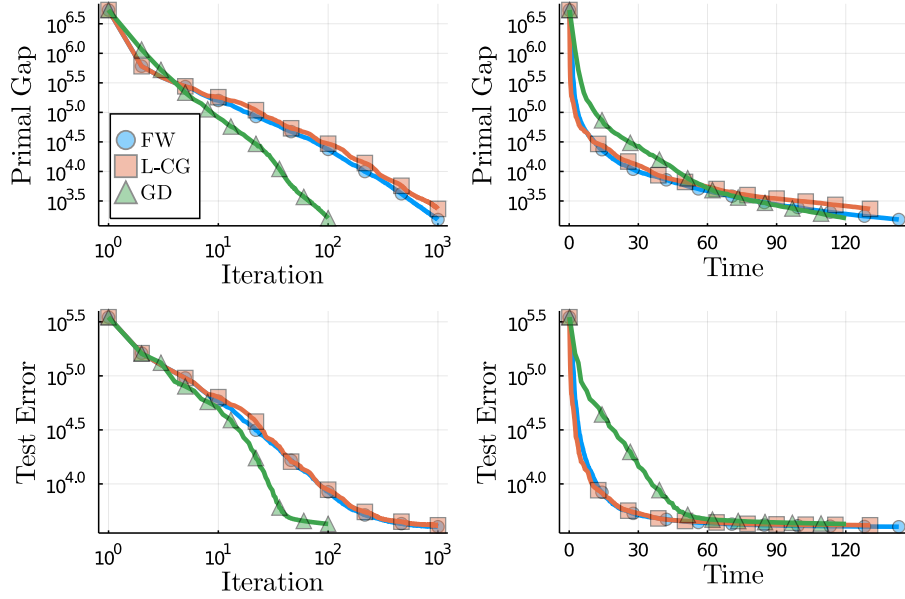


Figure 4: Comparison of standard and lazified FW with projected gradient descent on the Movielens example.

The results are presented in Fig. 4. We can clearly observe that the computational cost of a single PGD iteration is much higher than the cost of a FW variant step. The FW variants tested complete 10^3 iterations in around 120 seconds, while the PGD algorithm only completes 10^2 iterations in a similar time frame. We also observe that the progress per iteration made by each projection-free variant is smaller than the progress made by PGD, as expected. Note that, minimizing a linear function over the nuclear norm ball, in order to compute the LMO, amounts to computing the left and right singular vectors associated with the largest singular value, which we do using the ARPACK (Lehoucq et al., 1998) Julia wrapper in the current example. On the other hand, projecting onto the nuclear norm ball requires computing a full singular value decomposition. The underlying linear solver can be switched by users developing their own LMO.

The top two figures in Fig. 4 present the primal gap of Problem (5) in terms of iteration count and wall-clock time. The two bottom figures show the performance on a test set of entries. Note that the test error stagnates for all methods, as expected. Even though the training error decreases linearly for PGD for all iterations, the test error stagnates quickly. The final test error of PGD is about 6% higher than the final test error of the standard FW algorithm, which is also 2% smaller than the final test error of the lazy FW algorithm. We would like to stress though that the intention here is primarily to showcase the algorithms and the results are considered to be illustrative in nature only rather than a proper evaluation with correct hyper-parameter tuning.

Another key aspect of FW algorithms is the sparsity of the provided solutions. Sparsity in this context refers to a matrix being low-rank. Although each solution is a dense matrix in terms of non-zeros, it can be decomposed as a sum of a small number of rank-one terms, each represented as a pair of left and right vectors. At each iteration, FW algorithms add at most one rank-one term to the iterate, thus resulting in a low-rank solution by design. In our example here, the final FW solution is of rank at most 95 while the lazified version provides a sparser solution of rank at most 80. The lower rank of the lazified FW is due to the fact that this algorithm sometimes avoids calling the LMO if there already exists an atom (here rank-1 factor) in the cache that guarantees enough progress; the higher sparsity might help with interpretability and robustness to noise. In contrast, the solution computed by PGD is of full column rank and even after truncating the spectrum, removing factors with small singular values, it is still of much higher rank than the FW solutions.

4.3 Exact optimization with rational arithmetic

The package allows for exact optimization with rational arithmetic. For this, it suffices to set up the LMO to be rational and choose an appropriate step-size rule as detailed below. For the LMOs included in the package, this simply means initializing the radius with a rational-compatible element type, e.g., 1, rather than a floating-point number, e.g., 1.0. Given that numerators and denominators can become quite large in rational arithmetic, it is strongly advised to base the used rationals on extended-precision integer types such as `BigInt`, i.e., we use `Rational{BigInt}`. For the probability simplex LMO with a rational radius of 1, the LMO would be created as follows:

```
lmo = FrankWolfe.ProbabilitySimplexOracle{Rational{BigInt}}(1)
```

As mentioned before, the second requirement ensuring that the computation runs in rational arithmetic is a rational-compatible step-size rule. The most basic step-size rule compatible with rational optimization is the `agnostic` step-size rule with $\gamma_t = 2/(2+t)$. With this step-size rule, the gradient does not even need to be rational as long as the atom computed by the LMO is of a rational type. Assuming these requirements are met, all iterates and the computed solution will then be rational:

```
n = 100
x = fill(big(1)//100, n)
# equivalent to { 1/100 }^100
```

Another possible step-size rule is `rationalshortstep` which computes the step size by minimizing the smoothness inequality as $\gamma_t = \frac{\langle \nabla f(\mathbf{x}_t), \mathbf{x}_t - \mathbf{v}_t \rangle}{2L \|\mathbf{x}_t - \mathbf{v}_t\|^2}$. However, as this step size depends on an upper bound on the Lipschitz constant L as well as the inner product with the gradient $\nabla f(\mathbf{x}_t)$, both have to be of a rational type.

4.4 Formulating the LMO with MathOptInterface

In this example, we project a random point onto a ℓ_1 -norm ball with the basic Frank-Wolfe algorithm using either the specialized LMO defined in the package or a generic LP formulation using `MathOptInterface.jl` and `GLPK` as underlying LP solver.

```

using LinearAlgebra

import MathOptInterface
const MOI = MathOptInterface
import GLPK
using FrankWolfe

n = 1000
xp = rand(n)

f(x) = norm(x - xp)^2
function grad!(storage, x)
    @. storage = 2 * (x - xp)
    return nothing
end

lmo_radius = 2.5
# create a L1-norm ball of radius 2.5
lmo = FrankWolfe.LpNormLMO{Float64,1}(lmo_radius)

x0 = FrankWolfe.compute_extreme_point(lmo, zeros(n))
gradient = similar(x0)

x_lmo, _ = frank_wolfe(f, grad!, lmo, x0)

# create a MathOptInterface Optimizer and build the equivalent region
o = GLPK.Optimizer()
x = MOI.add_variables(o, n)

# add the constraints to the GLPK optimizer
# x_i ≥ 0
# sum(x_i) == 1

# define LMO from the GLPK instance
lmo_moi = FrankWolfe.MathOptLMO(o)

x_moi, _ = frank_wolfe(f, grad!, lmo_moi, x0)

```

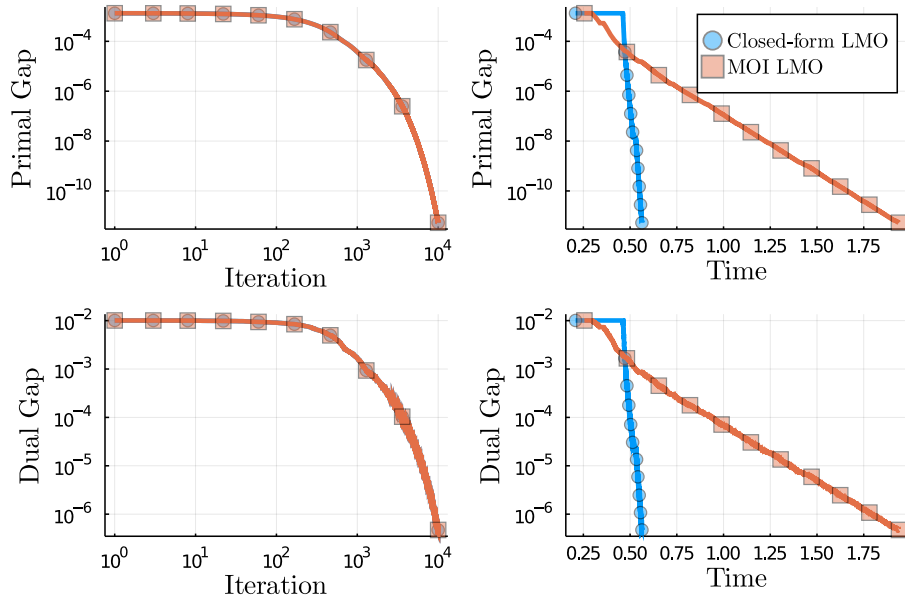


Figure 5: Performance comparison of custom LMO and MOI LP definition.

The resulting primal and dual progress are presented in Fig. 5. Since we use identical algorithmic parameters, the per-iteration progress is identical for the two LMOs. The closed-form LMO is however much faster to solve each iteration since it allocates only the output vector and avoids solving a generic linear optimization problem.

4.5 Doubly stochastic matrices

The set of doubly stochastic matrices or Birkhoff polytope appears in various combinatorial problems including matching and ranking. It is the convex hull of permutation matrices, a property of interest for FW algorithms because the individual atoms returned by the LMO only have n non-zero entries for $n \times n$ matrices. A linear function can be minimized over the Birkhoff polytope using the Hungarian algorithm. This LMO is substantially more expensive than minimizing a linear function over the ℓ_1 -ball norm, and thus the algorithm performance benefits from lazification. We present the performance profile of several FW variants in the following example on 200×200 matrices. The results are presented in Fig. 6.

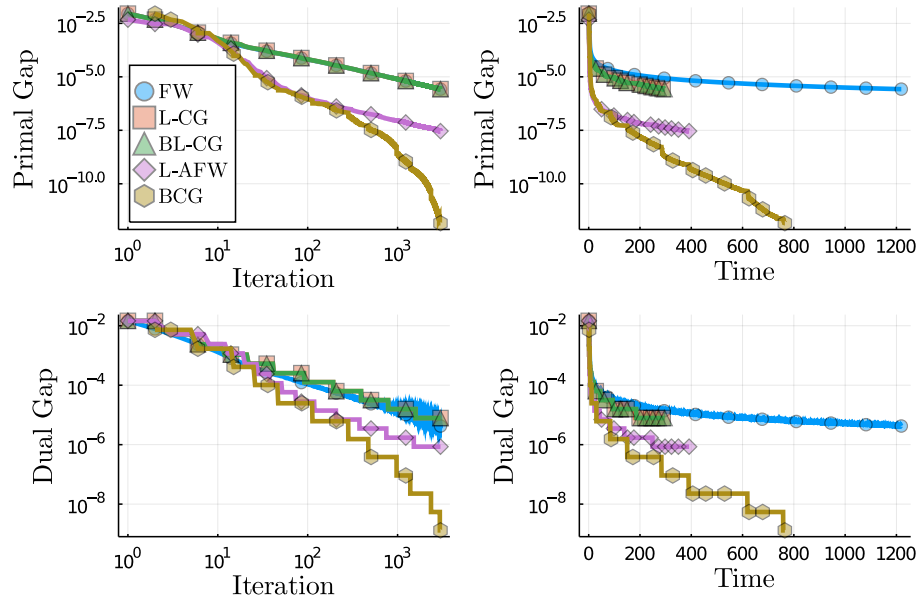


Figure 6: Lazified and eager variants on the Birkhoff polytope. FW stands for the classic Frank-Wolfe algorithm, L-CG uses a caching LMO with an unbounded cache size while BL-CG uses a bounded cache with 500 elements, L-AFW is the lazified AFW algorithm, and BCG is the Blended Conditional Gradient algorithm.

The per-iteration primal value evolution is nearly identical for FW and the lazy cache variants. We can observe a slower decrease rate in the first 10 iterations of BCG for both the primal value and the dual gap. This initial overhead is however compensated after the first iterations, BCG is the only algorithm terminating with the desired dual gap of 10^{-7} and not with the iteration limit. In terms of runtime, all lazified variants outperform the standard FW, the overhead of allocating and managing the cache are compensated by the reduced number of calls to the LMO.

5. Final Comments and Future Work

The `FrankWolfe.jl` package will be further extended over time and we welcome contributions, reporting of issues and bugs, as well as pull requests under the package [GitHub repository](#). A few prominent features to come in the near future will be:

1. Stronger interfacing with the broader Julia optimization ecosystem.
2. Interfacing with Python similar to, e.g., [diffeqpy](#).
3. Implementating closely related variants, such as, e.g., matching pursuit variants.

Acknowledgements

Research reported in this paper was partially supported through the Research Campus Modal funded by the German Federal Ministry of Education and Research (fund numbers 05M14ZAM, 05M20ZBM) and the Deutsche Forschungsgemeinschaft (DFG) through the DFG Cluster of Excellence MATH+.

References

- Antonello, N., Stella, L., Patrinos, P., and van Waterschoot, T. Proximal gradient algorithms: Applications in signal processing. *arXiv preprint arXiv:1803.01621*, 2018.
- Bertsimas, D., Cory-Wright, R., and Pauphilet, J. Mixed-projection conic optimization: A new paradigm for modeling rank constraints. *arXiv preprint arXiv:2009.10395*, 2020.
- Bezanson, J., Edelman, A., Karpinski, S., and Shah, V. B. Julia: A fresh approach to numerical computing. *SIAM review*, 59(1):65–98, 2017.
- Braun, G., Pokutta, S., and Zink, D. Lazifying conditional gradient algorithms. In *Proceedings of the 34th International Conference on Machine Learning*, pp. 566–575, 2017.
- Braun, G., Pokutta, S., Tu, D., and Wright, S. Blended conditional gradients: the unconditioning of conditional gradients. In *Proceedings of the 36th International Conference on Machine Learning*, 2019.
- Candès, E. J. and Plan, Y. Matrix completion with noise. *Proceedings of the IEEE*, 98(6):925–936, 2010.
- Candès, E. J. and Recht, B. Exact matrix completion via convex optimization. *Foundations of Computational mathematics*, 9(6):717–772, 2009.
- Candès, E. J. and Tao, T. The power of convex relaxation: Near-optimal matrix completion. *IEEE Transactions on Information Theory*, 56(5):2053–2080, 2010.
- Candès, E. J. and Wakin, M. B. An introduction to compressive sampling. *IEEE signal processing magazine*, 25(2):21–30, 2008.
- Candes, E. J., Romberg, J. K., and Tao, T. Stable signal recovery from incomplete and inaccurate measurements. *Communications on Pure and Applied Mathematics: A Journal Issued by the Courant Institute of Mathematical Sciences*, 59(8):1207–1223, 2006.
- Combettes, C. W. and Pokutta, S. Complexity of linear minimization and projection on some sets. *arXiv preprint arXiv:2101.10040*, 2021.
- Donoho, D. L. Compressed sensing. *IEEE Transactions on information theory*, 52(4):1289–1306, 2006.
- Dunning, I., Huchette, J., and Lubin, M. JuMP: A modeling language for mathematical optimization. *SIAM review*, 59(2):295–320, 2017.
- Fazel, M. *Matrix rank minimization with applications*. PhD thesis, PhD thesis, Stanford University, 2002.
- Frank, M. and Wolfe, P. An algorithm for quadratic programming. *Naval research logistics quarterly*, 3(1-2): 95–110, 1956.

- Guélat, J. and Marcotte, P. Some comments on Wolfe’s ‘away step’. *Mathematical Programming*, 35(1): 110–119, 1986.
- Hazan, E. and Luo, H. Variance-reduced and projection-free stochastic optimization. In *Proceedings of the 33th International Conference on Machine Learning*, pp. 1263–1271. PMLR, 2016.
- Juditsky, A. and Nemirovski, A. *Statistical Inference via Convex Optimization*, volume 69. Princeton University Press, 2020.
- Lacoste-Julien, S. and Jaggi, M. On the global linear convergence of Frank-Wolfe optimization variants. In *Proceedings of the 29th Conference on Neural Information Processing Systems*, pp. 566–575, 2015.
- Legat, B., Dowson, O., Garcia, J. D., and Lubin, M. MathOptInterface: a data structure for mathematical optimization problems. *arXiv preprint arXiv:2002.03447*, 2020.
- Legat, B., Timme, S., Deits, R., de Laat, D., Huchette, J., Saba, E., Forets, M., and Breiding, P. JuliaAlgebra/MultivariatePolynomials.jl: v0.3.13, April 2021. URL <https://doi.org/10.5281/zenodo.4656033>.
- Lehoucq, R. B., Sorensen, D. C., and Yang, C. *ARPACK users’ guide: solution of large-scale eigenvalue problems with implicitly restarted Arnoldi methods*. SIAM, 1998.
- Levitin, E. S. and Polyak, B. T. Constrained minimization methods. *USSR Computational Mathematics and Mathematical Physics*, 6(5):1–50, 1966.
- Mogensen, P. K. and Riseth, A. N. Optim: A mathematical optimization package for Julia. *Journal of Open Source Software*, 3(24):615, 2018. doi: 10.21105/joss.00615.
- Orban, D., Siqueira, A. S., and contributors. JSOSolvers.jl: JuliaSmoothOptimizers optimization solvers. <https://github.com/JuliaSmoothOptimizers/JSOSolvers.jl>, August 2020.
- Pedregosa, F., Negiar, G., Askari, A., and Jaggi, M. Linearly convergent Frank–Wolfe with backtracking line-search. In *Proceedings of the 23rd International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2020.
- Robbins, H. and Monro, S. A stochastic approximation method. *The Annals of Mathematical Statistics*, 22(3):400–407, 1951.
- Tropp, J. A. Just relax: Convex programming methods for identifying sparse signals in noise. *IEEE transactions on information theory*, 52(3):1030–1051, 2006.
- Udell, M. and Townsend, A. Why are big data matrices approximately low rank? *SIAM Journal on Mathematics of Data Science*, 1(1):144–160, 2019.
- Udell, M., Mohan, K., Zeng, D., Hong, J., Diamond, S., and Boyd, S. Convex optimization in Julia. In *2014 First Workshop for High Performance Technical Computing in Dynamic Languages*, pp. 18–28. IEEE, 2014.

Appendix A. Full nuclear norm example

Setup and data preprocessing:

```
# constructing the rating matrix from dataset
rating_matrix = sparse(
    ratings_frame[:, :userId],
    movies_indices,
    ratings_frame[:, :rating],
    length(users),
    length(movies),
)

# split present indices over test and training set
missing_rate = 0.05
missing_ratings = Tuple{Int,Int}[]
present_ratings = Tuple{Int,Int}[]
(I, J, V) = SparseArrays.findnz(rating_matrix)
for idx in eachindex(I)
    if V[idx] > 0
        if rand() <= missing_rate
            push!(missing_ratings, (I[idx], J[idx]))
        else
            push!(present_ratings, (I[idx], J[idx]))
        end
    end
end

function f(X)
    return 0.5 * sum(present_ratings) do (i,j)
        (X[i, j] - rating_matrix[i, j])^2
    end
end

function grad!(storage, X)
    storage .= 0
    for (i, j) in present_ratings
        storage[i, j] = X[i, j] - rating_matrix[i, j]
    end
    return nothing
end

function test_loss(X)
    return 0.5 * sum(missing_ratings) do (i,j)
        (X[i, j] - rating_matrix[i, j])^2
    end
end

# compute projection of X with nuclear norm of given radius
function project_nuclear_norm_ball(X; radius=1.0)
    U, sing_val, Vt = svd(X)
    if (sum(sing_val) <= radius)
        return (X, -norm_estimation * U[:, 1] * Vt[:, 1]')
    end
    σ = FrankWolfe.projection_simplex_sort(sing_val, s=radius)
    return U * Diagonal(σ) * Vt', -norm_estimation * U[:, 1] * Vt[:, 1]'
end

norm_estimation = 400 * Arpack.svds(
    rating_matrix, nsv=1, ritzvec=false,
)[1].S[1]

lmo = FrankWolfe.NuclearNormLMO(norm_estimation)
x0 = FrankWolfe.compute_extreme_point(lmo, zero(rating_matrix))
k = 100

gradient = spzeros(size(x0)...)
gradient_aux = spzeros(size(x0)...)

```

PGD and FW run:

```
xgd = Matrix(x0)
function_values = Float64[]
timing_values = Float64[]
function_test_values = Float64[]

# PGD steps
time_start = time_ns()
for _ in 1:k
    f_val = f(xgd)
    push!(function_values, f_val)
    push!(function_test_values, test_loss(xgd))
    push!(timing_values, (time_ns() - time_start) / 1e9)
    @info f_val
    grad!(gradient, xgd)
    xgd_new, vertex = project_nuclear_norm_ball(
        xgd - gradient / L_estimate, radius=norm_estimation,
    )
    gamma, _ = FrankWolfe.backtrackingLS(
        f, gradient, xgd, xgd - xgd_new, 1.0
    )
    @. xgd -= gamma * (xgd - xgd_new)
end

# pushes to the trajectory array the first 5 elements of the
# state and the test value at the current iterate
function build_callback(trajectory_arr)
    function callback(state)
        push!(
            trajectory_arr,
            Tuple(state)[1:5]..., test_loss(state.x)
        )
    end
end

trajectory_arr_fw = []
callback = build_callback(trajectory_arr_fw)
xfin, _, _, _, traj_data = FrankWolfe.frank_wolfe(
    f,
    grad!,
    lmo,
    x0;
    epsilon=1e-9,
    max_iteration=k,
    print_iter=k / 10,
    verbose=true,
    linesearch_tol=1e-8,
    line_search=FrankWolfe.Backtracking(),
    emphasis=FrankWolfe.memory,
    gradient=gradient,
    callback=callback,
)

trajectory_arr_lazy = []
callback = build_callback(trajectory_arr_lazy)
xlazy, _, _, _, _ = FrankWolfe.lazified_conditional_gradient(
    f,
    grad!,
    lmo,
    x0;
    epsilon=1e-9,
    max_iteration=k,
    print_iter=k / 10,
    verbose=true,
    linesearch_tol=1e-8,
    line_search=FrankWolfe.Backtracking(),
    emphasis=FrankWolfe.memory,
    gradient=gradient,
    callback=callback,
)
```