# Course Notes for EE227C (Spring 2018): Convex Optimization and Approximation

Instructor: Moritz Hardt

Email: `hardt+ee227c@berkeley.edu`

Graduate Instructor: Max Simchowitz

Email: `msimchow+ee227c@berkeley.edu`

January 31, 2018

## 5 Lecture 5: Conditional Gradient (Frank-Wolfe)

In this lecture we discussed about the conditional gradient method, also known as the Frank-Wolfe (FW) algorithm. The motivation of using this approach is the projected gradient descent can be computationally inefficient under certain scenarios.

### 5.1 Intuition behind Frank-Wolfe algorithm

We motivate this lecture by the computational inefficiency that projected gradient can have. With conditional gradient, we are able to sidestep some of these inefficiencies. An intuitive idea of the FW algorithm is as follows.

We start from $x_0$. Then, for t = 1 to T steps, we set

$$x_{t+1} = x_t + \eta_t(\bar{x}_t - x_t)$$

where

$$\bar{x}_t = \operatorname*{argmin}_{x \in \Omega} f(x_t) + \nabla f(x_t)^\top (x - x_t)$$

Since we hope to minimize with respect to $x \in \Omega$, we can simplify the equation.

$$\bar{x}_t = \operatorname*{argmin}_{x \in \Omega} \nabla f(x_t)^\top x$$

Note that we need step size $\eta_t \in [0, 1]$ to guarantee $x_{t+1} \in \Omega$.

## 5.2 Theorem: conditional gradient convergence analysis

**Theorem 5.1** (Convergence Analysis). *Assume:*

- *$f \colon \Omega \to \mathbb{R}$ is convex and $\beta$-smooth*

- *$\exists$ an optimal point, $x^*$, such that $x^* \in \Omega$ and $\nabla f(x^*) = 0$*

*Then, Frank-Wolfe achieves*

$$f(x_t) - f(x^*) \leqslant \frac{2\beta D^2}{t+2}$$

*with step size*

$$\eta_t = \frac{2}{t+2}$$

*where D is the diameter of $\Omega$, defined:*

$$D = \max_{x-y \in \Omega} \|x - y\|$$

Note that we can trade our assumption of the existence of $x^*$ for a dependence on L in our bound.

*Proof of Theorem 5.1.* By smoothness and convexity, we have

$$f(y) \leqslant f(x) + \nabla f(x)^\top (x - x_t) + \frac{\beta}{2}\|x - y\|^2$$

Letting $y = x_{t+1}$ and $x = x_t$, combined with the progress rule of conditional gradient descent, the above equation yields:

$$f(x_{t+1}) \leqslant f(x_t) + \eta_t \nabla f(x_t)^\top (\bar{x}_t - x_t) + \frac{\eta_t^2 \beta}{2}\|\bar{x}_t - x_t\|^2$$

We now recall the definition of $D$ from Theorem 5.1 and observe that $\|\bar{x}_t - x_t\|^2 \leqslant D^2$. Thus, we rewrite the inequality:

$$f(x_{t+1}) \leqslant f(x_t) + \eta_t \nabla f(x_t)^\top (x_t^* - x_t) + \frac{\eta_t^2 \beta D^2}{2}$$

Because of convexity, we also have that

$$\nabla f(x_t)^\top (x^* - x_t) \leqslant f(x^*) - f(x_t)$$

Thus,

$$f(x_{t+1}) - f(x^*) \leqslant (1 - \eta_t)(f(x_t) - f(x^*)) + \frac{\eta_t^2 \beta D^2}{2} \tag{1}$$

We use induction in order to prove $f(x_t) - f(x^*) \leqslant \frac{2\beta D^2}{t+2}$ based on Equation 1 above.

First step: Base Case $t = 0$

Since $f(x_{t+1}) - f(x^*) \leqslant (1 - \eta_t)(f(x_t) - f(x^*)) + \frac{\eta_t^2 \beta D^2}{2}$, when t=0, $\eta_t = \frac{2}{0+2} = 1, then$

$$f(x_1) - f(x^*) \leqslant (1 - \eta_t)(f(x_t) - f(x^*)) + \frac{\beta}{2} \|x_1 - x^*\|^2$$
$$= (1 - 1)(f(x_t) - f(x^*)) + \frac{\beta}{2} \|x_1 - x^*\|^2$$
$$\leqslant \frac{\beta D^2}{2}$$
$$\leqslant \frac{2\beta D^2}{3}$$

Thus, the induction hypothesis holds for our base case.

(Note that we cannot directly use $f(x_{t+1}) - f(x^*) \leqslant (1 - \eta_t)(f(x_t) - f(x^*)) + \frac{\eta_t^2 \beta D^2}{2}$ directly to derivate the base case out. Because $(1 - \eta_t)(f(x_t) - f(x^*))$ is greater than 0 and cannot be directly proved to be bounded.)

Second step: We assume that $f(x_t) - f(x^*) \leqslant \frac{2\beta D^2}{t+2}$ holds $\forall t$ and can show that the induction hyposthesis holds in general.

General case: given $f(x_t)$ holds for our destinated inequality, we'd like to show it also holds for $f(x_{t+1})$.

Using [Equation 1](),

$$f(x_{t+1}) - f(x) \leqslant (1 - \frac{2}{t+2})(f(x_t) - f(x^*)) + \frac{4}{2(t+2)}\beta D^2$$
$$\leqslant (1 - \frac{2}{t+2})(\frac{2\beta D^2}{t+2}) + \frac{4}{2(t+2)}\beta D^2$$
$$= \beta D^2(\frac{2t}{(t+2)^2} + \frac{2}{(t+2)^2})$$
$$= 2\beta D^2 \frac{t+1}{(t+2)^2}$$
$$= 2\beta D^2(\frac{t+1}{t+2})(\frac{1}{t+2})$$
$$\leqslant 2\beta D^2(\frac{t+2}{t+3})(\frac{1}{t+2})$$
$$= 2\beta D^2 \frac{1}{t+3}$$

Thus, the inequality also holds for the t+1 case.

■

## 5.3 Examples

The code for the following examples can be found [here]().

### 5.3.1 Nuclear norm projection

The *nuclear norm* (sometimes called *Schatten* 1-*norm* or *trace norm*) of a matrix $A$, denoted $\|A\|_*$, is defined as the sum of its singular values

$$\|A\|_* = \sum_i \sigma_i(A).$$

The norm can be computed from the singular value decomposition of $A$. We denote the unit ball of the nuclear norm by

$$B_*^{m \times n} = \{A \in \mathbb{R}^{m \times n} \mid \|A\|_* \leqslant 1\}.$$

How can we project a matrix $A$ onto $B_*$? Formally, we want to solve

$$\min_{X \in B_*} \|A - X\|_F^2$$

Due to the rotational invariance of the Frobenius norm, the solution is obtained by projecting the singular values onto the unit simplex. This operation corresponds to shifting all singular values by the same parameter $\theta$ and clipping values at $0$ so that the sum of the shifted and clipped values is equal to 1. This algorithm can be found from [DSSSC08].

### 5.3.2 Low-rank matrix completion

Suppose we have a partially observable matrix $Y$ and we would like to find its completion form projected in a nuclear norm ball. Formally we have the objective function

$$\min_{X \in B_*} \frac{1}{2} \|Y - X \odot O\|_F^2$$

And calculate the gradient of this function we will have

$$\nabla f(X) = Y - X \odot O$$

We can use projected gradient descent to solve this problem but it is more efficient to use Frank-Wolfe algorithm. We need to solve the linear optimization oracle

$$\bar{X}_t \in \operatorname*{argmin}_{X \in B_*} f(X_t) + \nabla f(X_t)^\top (X - X_t)$$

This will lead to a rank-1 matrix which is decomposed from $-\nabla f(X_t)$. This can be derived from Lemma 5.2. Now we can have the update rule for the conditional gradient as

$$X_{t+1} = X_t + \eta_t(-u_1 v_1^\top - X_t)$$

where $u_1$ and $v_1$ are the top left and right singular vectors.

**Lemma 5.2.**

$$\min_{X \in B_*} \langle \nabla f(X_t)^\top, X \rangle, \quad B_* = \{X | \|X\|_* \leqslant 1\}$$

*The optimal result is*

$$-\|\nabla f(X_t)\|$$

*with $X = -u_1 v_1^\top$ being one of the X that minimize the function where $u_1$ and $v_1$ are left and right singular vectors corresponding to the max singular value, $\|\nabla f(X_t)\| = \max_{\|X\|_* \leqslant 1} < \nabla f(X_t), X >$.*

*Proof of Lemma 5.2.* Since the dual norm of a nuclear norm is operator norm,

$$\|Y\| = \max_{\|X\|_* \leqslant 1} \langle Y, X \rangle$$

we can easily see that the optimal result of the objective function in Lemma 5.2 being

$$-\|Y\| = -\nabla f(X_t) = -\sigma_{\max}.$$

Then we show by Singular value decompasition to prove that $X = -u_1 v_1^\top$ is one of the solutions that minimize the function. From SVD, $\nabla f(X_t) = U\Sigma V^\top$, then $\Sigma = U^\top Y V$, we set $X = -u_1 v_1^\top$, thus

$$\langle \nabla f(X_t), X \rangle = \nabla f(X_t)^\top X = -V\Sigma U^\top u_1 v_1^\top = -\sigma_1$$

∎

According to the lemma, we can see that only the leading singular value and corresponding vectors are used to deprive the argmin X. Thus we can do the rank-1 approximation of this function which makes no difference to final result.

Power method works well in approximate the leading singular value and corresponding left and right vetors.
Here is the introduction of the process of SVD power method.

- $w_k = Av_{k-1}, \alpha_k = \|w_k\|, u_k = \alpha_k^{-1} w_k$

- $z_k = A^\top, \beta_k = \|z_k\|, v_k = \beta_k z_k$

$k = 1, 2, 3, 4, 5, ..., \sigma_{max} = \beta_k, v_k = \delta_{2k} \sum_{i=1}^{r} \sigma_j^{2k} y_j v_j, u_k = \delta_{2k+1} \sum_{i=1}^{r} \sigma_j^{2k+1} y_j v_j$

The basic idea behind is to use singular value. The proof of convergence can be found in many papers and sources. Here we provide a reference [BK15] and one can find the convergence analysis on pg 49.

# References

[BK15]     A. H. Bentbib and A. Kanber. Block Power Method for SVD Decomposition. *Analele Universitatii "Ovidius" Constanta - Seria Matematica*, 23(2), jan 2015.

[DSSSC08]  John Duchi, Shai Shalev-Shwartz, Yoram Singer, and Tushar Chandra. Efficient projections onto the l 1 -ball for learning in high dimensions. In *Proceedings of the 25th international conference on Machine learning - ICML '08*, pages 272–279, New York, New York, USA, 2008. ACM Press.