



STOCHASTIC OPTIMIZATION AND AUTOMATIC
DIFFERENTIATION FOR MACHINE LEARNING

Etude de l'algorithme Stochastic Dual Coordinate Ascent (SDCA)

Vincent LE MEUR

13 mai 2018

Table des matières

Introduction	2
Mise en place du SDCA	3
Mise en place de PEGASOS	4
Analyse	6
Conclusion	9

Introduction

L'objectif de cette étude est de résoudre un problème de Support Vector Machines (SVM) dans le cadre d'une classification binaire.

Cette étude se base sur l'article suivant : "[Stochastic Dual Coordinate Ascent Methods for Regularized Loss Minimization](#) (Shai Shalev-Shwartz, Tong Zhang)"

L'essentiel du code se trouve dans le notebook appelé SDCA-LE-MEUR.ipynb. Ce rapport comporte uniquement le code de définition des algorithmes principaux de l'étude ainsi que les graphs obtenus et leurs interprétations.

Nous allons utiliser les notations suivantes :

- $X_1, \dots, X_n \in \mathbb{R}^d$ pour les variables explicatives
- $y_1, \dots, y_n \in \{-1, 1\}$ pour les labels associés
- λ le paramètre de régularisation propre à l'algorithme
- $w \in \mathbb{R}$ le paramètre à optimiser du problème primal
- $\phi_i(a) = \max(0, 1 - y_i a)$ la fonction de perte (Hinge loss) associée au SVM étudié

Le problème primal associé à l'optimisation de notre SVM consiste alors avec les notations précédentes à minimiser l'application pénalisée suivante : $P(w) = \frac{1}{n} \sum_{i=1}^n \phi_i(w^T X_i) + \frac{\lambda}{2} \|w\|^2$ (1)

Une approche alternative (Dual Coordinate Ascent) à ce problème est celle qui étudie le problème dual. En prenant le conjugué convexe de notre fonction de perte ϕ précédente définie par $\phi_i^*(u) = \max_z (zu - \phi_i(z))$, nous pouvons en effet définir le problème équivalent au précédent qui est de maximiser : $D(\alpha) = \frac{1}{n} \sum_{i=1}^n -\phi_i^*(-\alpha_i) - \frac{\lambda}{2} \left\| \frac{1}{\lambda n} \sum_{i=1}^n \alpha_i x_i \right\|^2$ où ici $\alpha \in \mathbb{R}^n$ (2)

Nous voyons un lien clair entre les deux problèmes. En effet, soit l'application $\omega(\alpha) = \frac{1}{\lambda n} \sum_{i=1}^n \alpha_i x_i$ et α^* la solution du problème dual (2) alors $\omega^* = \omega(\alpha^*)$ est également la solution optimale du problème primal (1)

Nous allons alors étudier une approche utilisant cette dualité (SDCA) et la comparer à une approche de descente de sous Gradient :

- Stochastic Dual Coordinate Ascent (SDCA)
- Primal Estimated subGrAdient Solver for SVM (PEGASOS)

Mise en place du SDCA

L'algorithme SDCA consiste donc à résoudre de manière itérative les deux problèmes de la manière suivante :

Procedure SDCA($\alpha^{(0)}$)

Let $w^{(0)} = w(\alpha^{(0)})$

Iterate: for $t = 1, 2, \dots, T$:

Randomly pick i

Find $\Delta\alpha_i$ to maximize $-\phi_i^*(-(\alpha_i^{(t-1)} + \Delta\alpha_i)) - \frac{\lambda n}{2} \|w^{(t-1)} + (\lambda n)^{-1} \Delta\alpha_i x_i\|^2$

$\alpha^{(t)} \leftarrow \alpha^{(t-1)} + \Delta\alpha_i e_i$

$w^{(t)} \leftarrow w^{(t-1)} + (\lambda n)^{-1} \Delta\alpha_i x_i$

Output (Averaging option):

Let $\bar{\alpha} = \frac{1}{T-T_0} \sum_{i=T_0+1}^T \alpha^{(i-1)}$

Let $\bar{w} = w(\bar{\alpha}) = \frac{1}{T-T_0} \sum_{i=T_0+1}^T w^{(i-1)}$

return \bar{w}

Output (Random option):

Let $\bar{\alpha} = \alpha^{(t)}$ and $\bar{w} = w^{(t)}$ for some random $t \in T_0 + 1, \dots, T$

return \bar{w}

L'algorithme comporte deux types de sorties :

- Une sortie moyennée qui correspond à prendre pour ω final une moyenne des dernières itérations (ce nombre d'itérations considérées est fixé par le paramètre T_0)
- Une sortie "aléatoire" ou l'on sélectionne ω final aléatoirement parmi les dernières itérations (toujours fixées par le paramètre T_0)

On cherche à chaque étape à calculer $\Delta\alpha_i$. Il se trouve que pour le problème précis que l'on se fixe (résoudre un SVM avec perte Hinge) , on connaît explicitement cette quantité qui est : $\Delta\alpha_i = \max(-1, \min(\frac{y_i - x_i^T \omega^{t-1}}{\|x_i\|^2 / (\lambda n)} + \alpha_i^{t-1})) - \alpha_i^{t-1}$

Nous allons mettre en place deux variantes de cet algorithme : SDCA et SCDCAperm. La différence entre ces deux algorithmes est la manière d'itérer. Pour SDCA à chaque itération, on pioche uniformément avec remise un $i \in \{1, \dots, n\}$ alors que pour SCDCAperm on effectue à l'avance une permutation aléatoire de $\{1, \dots, n\}$ sur laquelle on travaille.

```

1 def SDCA(nb_epoch,T0,X,y,Lambda,method):
2     n=X.shape[0]
3     list_alpha=[]
4     list_omega=[]
5     list_primal=[]
6     list_dual=[]
7     alpha=np.zeros(n)
8     T=nb_epoch*n
9     omega = (1/(Lambda*n)) * (X.T.dot(alpha))
10    for t in range(1,T):
11        i=int(np.random.uniform(0,n))
12        delta_alpha=y[i]*max(0,min(1, (1-np.dot(X[i].T,np.dot(omega,y[i])))/(np.linalg.norm(X[
13            i]))/(Lambda*n)) +alpha[i]*y[i])) -alpha[i]
14        alpha[i]=alpha[i]+delta_alpha
15        list_alpha.append(alpha)
16        omega=omega + (1/(Lambda*n))*delta_alpha*X[i]
17        list_omega.append(omega)
18        list_primal.append(func_primal(omega,Lambda,n,X,y))
19        list_dual.append(func_dual(alpha,Lambda,n,X,y))
20    if method=="random":
21        k=int(np.random.uniform(T0,T))
22        omega_final=list_omega[k]
23        alpha_final=list_alpha[k]
24    if method=="average":
25        omega_final=(1/(T-T0))*sum(list_omega[k] for k in range(T0,T-1))
26        alpha_final=(1/(T-T0))*sum(list_alpha[k] for k in range(T0,T-1))
27    return(omega_final,list_primal,list_dual,list_alpha)

```

Mise en place du SDCA

```

1 def SDCA_perm(nb_epoch,T0,X,y,Lambda,method):
2     n=X.shape[0]
3     list_alpha=[]
4     list_omega=[]
5     list_primal=[]
6     list_dual=[]
7     alpha=np.zeros(n)
8     T=nb_epoch*n
9     omega = (1/(Lambda*n)) * (X.T.dot(alpha))
10    for k in range(1,nb_epoch):
11        # On fixe une permutation sur {1,...,n}
12        perm=permutation(n)
13        for i in range(len(perm)):
14            # On calcule explicitement le delta_alpha_i
15            delta_alpha=y[i]*max(0,min(1, (1-np.dot(X[i].T,np.dot(omega,y[i])))/(np.linalg.
16                norm(X[i]))/(Lambda*n)) +alpha[i]*y[i])) -alpha[i]
17            alpha[i]=alpha[i]+delta_alpha
18            list_alpha.append(alpha)
19            omega=omega + (1/(Lambda*n))*delta_alpha*X[i]
20            list_omega.append(omega)
21            list_primal.append(func_primal(omega,Lambda,n,X,y))
22            list_dual.append(func_dual(alpha,Lambda,n,X,y))
23    if method=="random":
24        k=int(np.random.uniform(T0,T))
25        omega_final=list_omega[k]
26        alpha_final=list_alpha[k]
27    if method=="average":
28        omega_final=(1/(T-T0))*sum(list_omega[k] for k in range(T0,T-1))
29        alpha_final=(1/(T-T0))*sum(list_alpha[k] for k in range(T0,T-1))
30    return(omega_final,list_primal,list_dual,list_alpha)

```

Mise en place du SDCAPerm

Mise en place de PEGASOS

Nous allons ici mettre en place l'algorithme "concurrent" Pegasos qui effectue une descente de gradient stochastique sur notre problème primal (1). Nous donnons ci-dessous son "pseudo code" :

```

INPUT:  $S, \lambda, T$ 
INITIALIZE: Set  $\mathbf{w}_1 = 0$ 
FOR  $t = 1, 2, \dots, T$ 
    Choose  $i_t \in \{1, \dots, |S|\}$  uniformly at random.
    Set  $\eta_t = \frac{1}{\lambda t}$ 
    If  $y_{i_t} \langle \mathbf{w}_t, \mathbf{x}_{i_t} \rangle < 1$ , then:
        Set  $\mathbf{w}_{t+1} \leftarrow (1 - \eta_t \lambda) \mathbf{w}_t + \eta_t y_{i_t} \mathbf{x}_{i_t}$ 
    Else (if  $y_{i_t} \langle \mathbf{w}_t, \mathbf{x}_{i_t} \rangle \geq 1$ ):
        Set  $\mathbf{w}_{t+1} \leftarrow (1 - \eta_t \lambda) \mathbf{w}_t$ 
    [ Optional:  $\mathbf{w}_{t+1} \leftarrow \min \left\{ 1, \frac{1/\sqrt{\lambda}}{\|\mathbf{w}_{t+1}\|} \right\} \mathbf{w}_{t+1}$  ]
OUTPUT:  $\mathbf{w}_{T+1}$ 

```

Ici encore nous allons donner deux variantes de cet algorithme. Une variante pegasos qui, comme indiqué ci-dessus prends un i à chaque itération uniformément parmi $\{1, \dots, n\}$, et une variante pegasosperm qui à chaque epoch fixe une permutation aléatoire de $\{1, \dots, n\}$ et itère sur celle-ci :

```

1 def pegasos(nb_epoch, X, y, Lambda):
2     n=X.shape[0]
3     d=X.shape[1]
4     list_omega=[]
5     list_primal=[]
6     omega =np.ones(d)
7     T=nb_epoch*n
8     for t in range(1,T):
9         i=int(np.random.uniform(0,n))
10        eta_t=1/(Lambda*t)
11        if np.dot(X[i].T,np.dot(omega,y[i]))<1:
12            omega=(1-eta_t*Lambda)*omega + eta_t*y[i]*(X[i])
13            list_omega.append(omega)
14            list_primal.append(func_primal(omega,Lambda,n,X,y))
15        else:
16            omega=(1-eta_t*Lambda)*omega
17            list_omega.append(omega)
18            list_primal.append(func_primal(omega,Lambda,n,X,y))
19    return(omega,list_primal)

```

Mise en place de PEGASOS

```

1 def pegasos_perm(nb_epoch,X,y,Lambda):
2     n=X.shape[0]
3     d=X.shape[1]
4     list_omega=[]
5     list_primal=[]
6     omega =np.ones(d)
7     T=nb_epoch*n
8     t=1
9     for k in range(1,nb_epoch):
10         perm=permutation(n)
11         for i in range(len(perm)):
12             eta_t=1/(Lambda*t)
13             if np.dot(X[i].T,np.dot(omega,y[i]))<1:
14                 omega=(1-eta_t*Lambda)*omega + eta_t*y[i]*(X[i])
15                 list_omega.append(omega)
16                 list_primal.append(func_primal(omega,Lambda,n,X,y))
17             else:
18                 omega=(1-eta_t*Lambda)*omega
19                 list_omega.append(omega)
20                 list_primal.append(func_primal(omega,Lambda,n,X,y))
21         t=t+1
22     return(omega,list_primal)

```

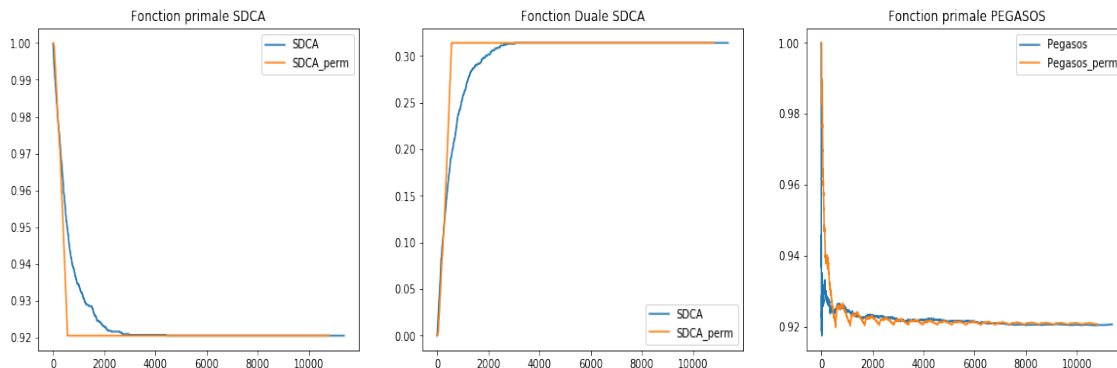
Mise en place de PEGASOSPERM

Analyse

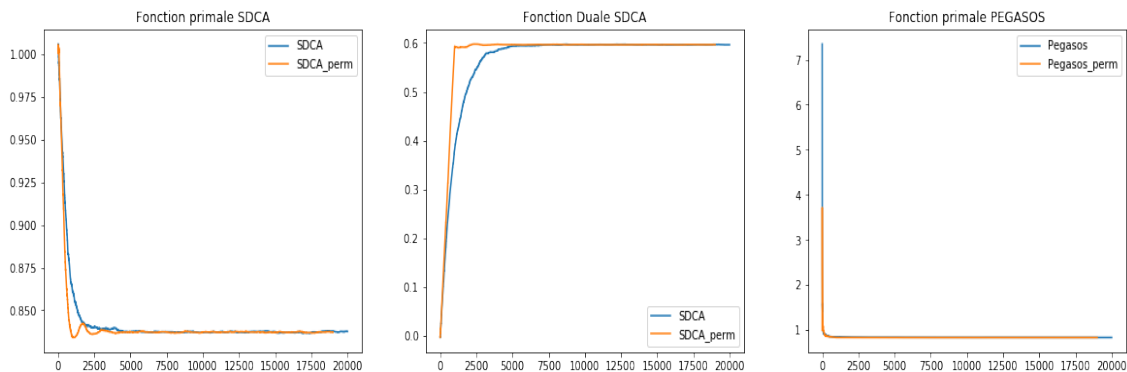
Nous allons à présent analyser nos algorithmes précédents sur deux datasets distincts. Nous allons regarder l'évolution de la convergence de notre minimisation ainsi que la précision finale obtenue dans le problème de classification.

Le premier correspond à un dataset de détection du cancer du sein disponible sur sklearn. Il comporte 569 observations pour 30 variables. Le second dataset est un dataset artificiel effectué à partir de réalisations bruitées d'une loi normale. Il comporte 1000 observations pour 50 variables.

Nous allons dans un premier temps évaluer pour chacun des algorithmes (SDCA et PEGASOS) la différence de convergence pour les deux variantes (avec ou sans permutation) pour les données du cancer du sein :



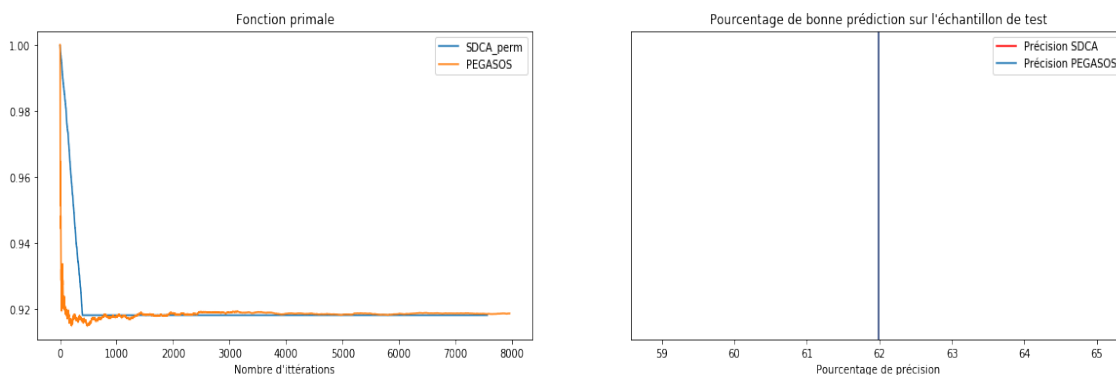
Nous effectuons la même étude concernant cette fois les données simulées :



En terme de convergence la variante "permutation" semble mieux "performer" sur les deux datasets pour l'algorithme SDCA, en revanche l'efficacité de la version "permutation" de PEGASOS a des performances plus discutables sur le premier dataset. Pour la suite nous garderons alors les algorithmes SDCAperm et pegasos pour notre étude.

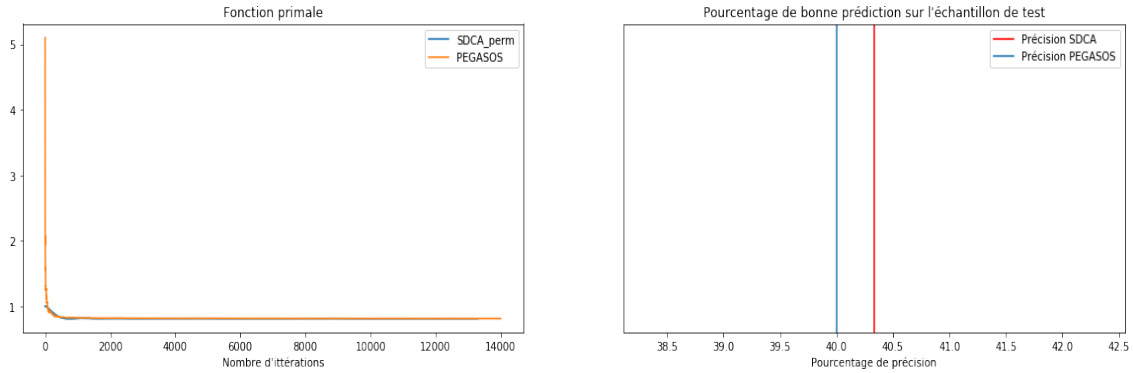
Jusqu'ici, nous avons considéré l'intégralité de nos données pour minimiser la fonction. Nous allons ici utiliser le découpage dataset d'apprentissage/dataset de test pour évaluer la précision des algorithmes également (cf code dans le notebook)

Commençons par le premier dataset portant sur le cancer du sein :



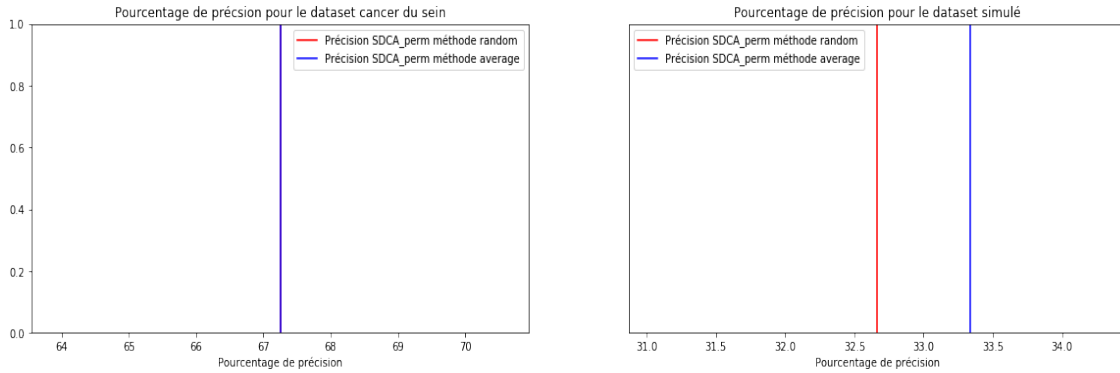
Nos deux algorithmes tendent quasiment vers le même ω ce qui implique un taux de précision similaire (environ 60 %) qui est plutôt faible. En terme de convergence, il semble que l'algorithme PEGASOS décroît plus rapidement vers le ω_{final} mais mets plus de temps à se stabiliser que l'algorithme SDCAperm. En effet l'algorithme PEGASOS semble se stabiliser au bout de 3000 itérations environ contre 500 pour le SDCA.

Effectuons à présent la même étude pour le second dataset :

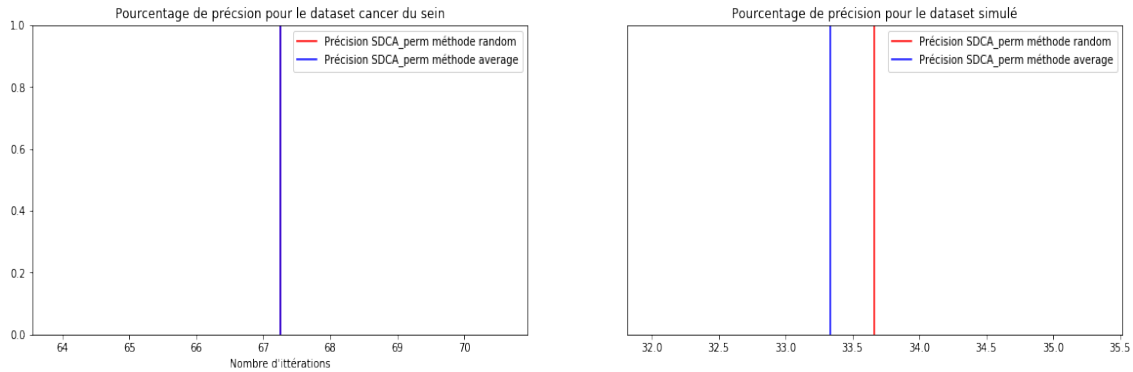


Sur ce dataset, les deux convergences sont très comparables en terme de vitesse et de stabilité. En revanche sur la précision, l'algorithme SDCA semble faire légèrement mieux que son concurrent PEGASOS. Dans les deux cas, puisque nous sommes dans un contexte de classification binaire les résultats ici sont plutôt faibles.

Nous allons à présent observer l'influence de la méthode choisie (random ou average) pour le SDCA sur la précision finale obtenue :



Nous constatons que pour le premier dataset les deux méthodes offrent une précision équivalente. Pour le second dataset un léger écart de précision apparaît. Ceci peut s'expliquer par la dimension plus importante (50 contre 30) du second dataset. La différence d'efficacité des deux méthodes est également liée au paramètre T_0 . Ce paramètre permet de fixer à partir de quelle itération nous allons tirer le ω_{final} aléatoirement ou à partir de quelle itération allons nous effectuer la moyenne selon la méthode employé. Etant donné que le comportement de l'algorithme est en phase transitoire sur les premières itérations, la méthode "average" risque d'être moins performante si T_0 est trop petit. A titre d'illustration nous traçons ci dessous la même étude avec un T_0 dix fois moins important :



On constate alors que la méthode "average" qui avait de meilleures performances que la méthode "random" sur le second dataset réussit moins bien ici. Ceci s'explique par l'influence du régime transitoire qui, si T_0 est trop petit, peut pénaliser cette méthode.

Conclusion

Nous avons dans notre étude implémenté l'algorithme SDCA ainsi qu'une variante avec une stratégie sous Gradient (PEGASOS). Nous avons alors comparé ces algorithmes sous deux variantes, en terme de minimisation du problème primal et de précision de la classification associé sur deux datasets distincts.

Les deux algorithmes ont quasiment les mêmes performances en terme de précision sur nos deux datasets (avec un léger avantage peut être pour le SDCA). Concernant la minimisation de la fonction primale notre étude semble montrer que l'algorithme PEGASOS descend plus vite vers un minimum de la fonction mais met plus de temps que SDCA à se stabiliser sur les premières itérations.