

## ▼ Computational statistics - TP 1

### LACOMBE Yoach

### Exercice 3

#### question 1

Il s'agit de faire une descente de gradient stochastique (SGD) pour minimiser le risque empirique. Ici le risque empirique est :

$$R_n(w) = \sum_{i=0}^{n-1} (y_i - w^T x_i)^2 / n = \sum_{i=0}^{n-1} j(w, z_i) / n$$

Décrivons rapidement l'algo de SGD:

1.  $w_0$  donné
2.  $k \rightarrow k+1, w_{k+1} = w_k - \epsilon_k \nabla_w j(w_k, z_{k+1})$

Il suffit donc maintenant de calculer le gradient de  $j$  par rapport à  $w$  et de décider du pas  $\epsilon_k$ .

On prend  $\epsilon_k = 1/k^\alpha$  avec  $\alpha$  dans  $]0.5, 1]$ .

Le gradient est très facile à calculer,  $\nabla_w j(w_k, z_{k+1}) = -2 * (y_{k+1} - w_k^T x_{k+1}) x_{k+1}$ .

Ainsi, l'étape d'update revient à :

$$\begin{aligned} w_{k+1} &= w_k + 2/k^\alpha * (y_{k+1} - w_k^T x_{k+1}) x_{k+1} \\ w_{k+1}^T &= w_k^T + 2/k^\alpha * (y_{k+1} - w_k^T x_{k+1}) x_{k+1}^T \end{aligned}$$

Ecrivons maintenant le code. Je choisis comme critère d'arrêt un nombre d'itérations fixe et un  $w_0$  aléatoire (tiré d'un vecteur gaussien centré réduit).

**Remarque 1:** On peut ajouter une dimension aux samples,  $z = [z^T, 1]^T$ , afin d'avoir le biais calculé dans l'algorithme.

**Remarque 2:** Je prends le paradigme de sklearn, c'est-à-dire un dataset de dimension  $(n, d)$ , càd  $data = [x_0, x_1, \dots, x_n]^T$

```
import numpy as np

def sgd(df, labels, alpha = 0.7, nb_iterations = 10000):
    df = np.c_[ df, np.ones(len(df))]

    n, d = df.shape

    w_transpose = np.random.multivariate_normal(np.zeros(d), np.identity(d))
    w_transpose = np.expand_dims(w_transpose, axis=0)
    for i in range(nb_iterations):
        j = i%n
```

```
w_transpose = w_transpose + 2*(labels[j] - w_transpose@df[j,:].T)*df[j,:]/((i+1)**alph

return w_transpose[0,:]
```

## ▼ question 2

Il faut décider de comment tirer les samples. Pour ça, je décide de tirer les samples d'une loi uniforme  $U([u_{min}, u_{max}]^d)$ .

```
import matplotlib.pyplot as plt
n = 2000
u_min = -1
u_max = 1
w_true = np.array([5,5, 1])

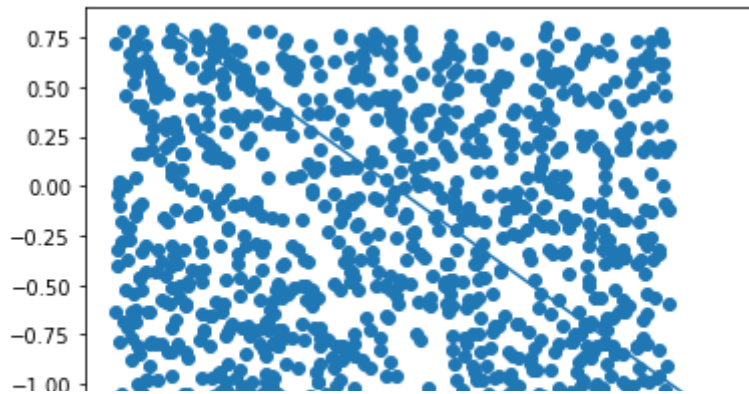
y_min = (-w_true[0]*u_min - w_true[2])/w_true[1]
y_max = (-w_true[0]*u_max - w_true[2])/w_true[1]
data = (y_max - y_min)*np.random.random_sample((n,2)) + y_min

plt.scatter(data[:1000,0], data[:1000,1])
plt.plot([u_min,u_max],[y_min, y_max ])
plt.show()

labels = np.zeros(n) - 1

labels[(data@w_true[:2] + w_true[2])>=0] = 1

plt.scatter(data[:1000,0], data[:1000,1], c = labels[:1000])
plt.plot([u_min,u_max], [(-w_true[0]*u_min - w_true[2])/w_true[1], (-w_true[0]*u_max - w_t
plt.show()
```



### question 3

```

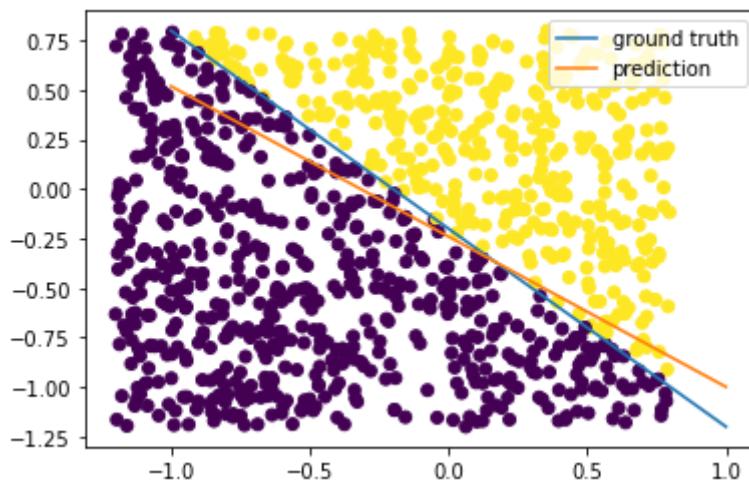
w_pred = sgd(data, labels, alpha = 0.5, nb_iterations=2000)

print("prediction :", w_pred)
print("true_value :", w_true)

prediction : [0.86639509 1.14158719 0.27417013]
true_value : [5 5 1]

plt.scatter(data[:1000,0], data[:1000,1], c = labels[:1000])
plt.plot([u_min,u_max], [(-w_true[0]*u_min - w_true[2])/w_true[1], (-w_true[0]*u_max - w_t
plt.plot([u_min,u_max], [(-w_pred[0]*u_min - w_pred[2])/w_pred[1], (-w_pred[0]*u_max - w_p
plt.legend()
plt.show()

```



Les hyperplans se ressemblent plus ou moins. La méthode paraît donc marcher. On remarque également que le vecteur prédit est plus ou moins le vrai vecteur à une constante multiplicative près.

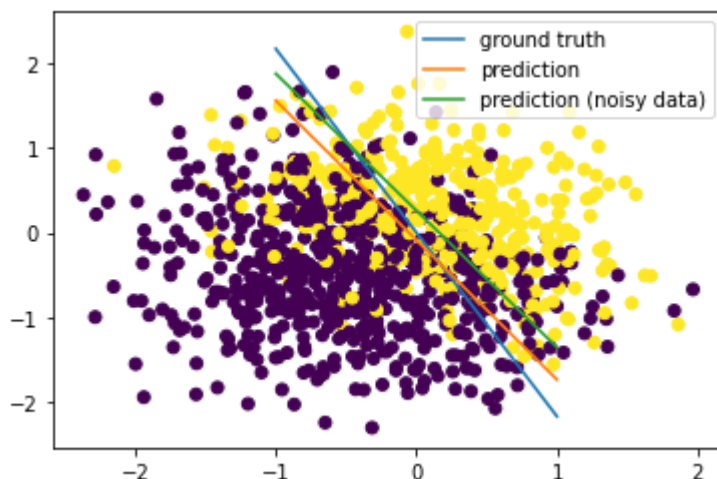
### question 4

```
noise = np.random.normal(0, 0.5, size = data.shape)
data_noise = data + noise
w_pred_noise = sgd(data_noise, labels, alpha = 0.5, nb_iterations=5000)
```

```
print("prediction :", w_pred)
print("prediction (noisy data):", w_pred_noise)
print("true_value :", w_true)
```

```
prediction : [0.86639509 1.14158719 0.27417013]
prediction (noisy data): [ 0.42469729  0.57047459 -0.0637899 ]
true_value : [5 5 1]
```

```
plt.scatter(data_noise[:1000,0], data_noise[:1000,1], c = labels[:1000])
plt.plot([u_min,u_max], [(-w_true[0]*np.min(data_noise[:1000,0]) - w_true[2])/w_true[1], (
plt.plot([u_min,u_max], [(-w_pred[0]*np.min(data_noise[:1000,0]) - w_pred[2])/w_pred[1], (
plt.plot([u_min,u_max], [(-w_pred_noise[0]*np.min(data_noise[:1000,0]) - w_pred_noise[2])/
plt.legend()
plt.show()
```



Ici la prédiction avec noisy data n'a pas l'air d'avoir marché (il faudrait regarder l'accuracy de la prédiction pour en être sûr).

**Remarque:** L'algorithme ne marche pas si les données ne sont pas normalisées (numériquement les valeurs explosent trop vite pour Python). Je normaliserai donc les données pour la prochaine question.

## ▼ question 5:

```
from sklearn.datasets import load_breast_cancer

X, y = load_breast_cancer(return_X_y = True)

from sklearn.preprocessing import normalize
```

```
X_normalized = normalize(X)
```

```
y[y==0] = -1
```

```
w_decision = sgd(X_normalized, y, alpha = 0.5, nb_iterations=len(X)*3)
```

Je teste l'algorithme avec l'accuracy (on suppose donc un bias égale à `w_decision[-1]` -> on aurait pu le choisir afin d'atteindre le trade-off souhaitée entre precision et recall, je ne fais pas ce travail ici).

```
y_pred = np.c_[ X_normalized, np.ones(len(X_normalized)) ]@w_decision  
y_pred = np.sign(y_pred)
```

```
print('accuracy', np.mean(y_pred == y))
```

```
accuracy 0.9068541300527241
```

---

✓ 0 s terminée à 15:12



### Exercice 2: (3)

Soit  $x \in \{\frac{1}{m}\}_{m \in \mathbb{N}^*}$ ,  $f$  mesurable bornée,

$$\begin{aligned}
 P^0(x) &= \mathbb{E}[f(X_1) | X_0 = x] \text{ avec la notation du TD} \\
 &= \mathbb{E}_0[f(U_0) | X_0 = x] \text{ (} U_0 \sim \mathcal{U}([0,1]) \perp X_0 \text{)} \\
 &= \mathbb{E}[f(U_0)] \text{ car } U_0 \perp X_0 \\
 &= \int f(t) \pi(t) dt \text{ avec } \pi \text{ la distribution de } \mathcal{U}([0,1]) \\
 &= \text{cte}
 \end{aligned}$$

Par récurrence, on montre facilement que  $\forall m \in \mathbb{N}^*$ ,  $P^m f(x) = \int f(t) \pi(t) dt$

car  $P^{m+1} f(x) = P^0[P^m f(x)] \stackrel{\text{rem}}{=} \mathbb{E}[P^m f(X_1) | X_0 = x]$

$$\begin{aligned}
 &= \mathbb{E}\left[\underbrace{\int f(t) \pi(t) dt}_{\text{cte}} | X_0 = x\right] \\
 &= \int f(t) \pi(t) dt
 \end{aligned}$$

(4) Soit  $m \geq 2$  et  $x = \frac{1}{m}$ .

a)  $m=2$ ,  $P^2(x, \frac{1}{m+2}) = \int P(y, \frac{1}{m+2}) P^0(x, dy)$

$\forall n$   $P^0(x, dy) = (1-x^2) \delta_{\frac{x}{1+x^2}}(dy)$  [car  $\int \frac{1}{1+t^2} dt = \arctan(t) \in [0, \pi/2]$  pour  $t \in [0,1]$ ]

D'où  $P^2(x, \frac{1}{m+2}) = P^0(1-x^2) \int P(y, \frac{1}{m+2}) \delta_{\frac{x}{1+x^2}}(dy)$

Je m'arrête par manque de définitions claires pour  $P^m$  (est-ce  $P^{m-2} \circ P$  ou  $P \circ P^{m-2}$  ?).

J'utilise donc cette déf:

$$P^{n+1}(x, A) = \int P(x, dy) P^n(y, A) \quad \text{et} \quad P^2(x, A) = P(x, A).$$

On a donc pour  $m=2$ ,  $P^2(\frac{1}{m}, \frac{1}{m+2}) = P(\frac{1}{m}, \frac{1}{m+2}) = 1 - \frac{1}{m^2}$ .

Par récurrence  $\boxed{\forall m \in \mathbb{N}^*, P^m(\frac{1}{m}, \frac{1}{m+m}) = \prod_{i=0}^{m-1} (1 - \frac{1}{(m+i)^2})}$

En effet  $P^{n+1}(\frac{1}{m}, \frac{1}{m+n+2}) = \int P(x, dy) P^n(y, \frac{1}{m+n+2}) = \int (1 - \frac{1}{m^2}) \times \delta_{\frac{x}{1+x^2}}(dy) P^n(y, \frac{1}{m+n+2})$

(récurrence)

$$\begin{aligned}
 &= (1 - \frac{1}{m^2}) P^n(\frac{1}{m+2}, \frac{1}{m+n+2}) = (1 - \frac{1}{m^2}) \times \prod_{i=0}^{n-1} (1 - \frac{1}{(m+2+i)^2}) = \prod_{i=0}^n (1 - \frac{1}{(m+i)^2}) \quad \square
 \end{aligned}$$

Q2) Remarquons que  $\pi(A) = \pi\left(\bigcup_{q \in \mathbb{N}} \left\{ \frac{1}{m+1+q} \right\}\right) = 0$  car l'ensemble est  
 de mesure de Lebesgue  
 nul.

Considérons  $v_m := p^m\left(\frac{x}{m}, \frac{1}{m+m}\right) = \prod_{i=0}^{m-1} \left(1 - \frac{1}{(m+i)^2}\right)$

$$\begin{aligned}
 \text{On a } v_m &= \prod_{i=0}^{m-1} \frac{(m+i-1)(m+i+1)}{(m+i)^2} = \underbrace{\prod_{i=0}^{m-1} \frac{m+i-1}{m+i}}_{= \frac{m-1}{m+m-1}} \times \underbrace{\prod_{i=0}^{m-1} \frac{m+i+1}{m+i}}_{= \frac{m+m}{m}} \\
 &= \frac{m-1}{m+m-1} \times \frac{m+m}{m}
 \end{aligned}$$

Donc  $v_m \sim \left(2 - \frac{1}{m}\right) \times \frac{m}{m}$

On a donc  $v_m \xrightarrow{m \rightarrow +\infty} 1 - \frac{1}{m} \quad (\neq 0 \text{ pour } m \geq 2)$

Puisque  $\forall m \in \mathbb{N}^*$ ,  $a_m := \frac{1}{m+m} \in A$

on a,  $p^m(x, A) \geq \underbrace{p^m(x, a_m)}_{\xrightarrow{m \rightarrow +\infty} 1 - \frac{1}{m} > 0}$

Ainsi  $p^m(x, A)$  ne CV pas vers 0 donc ne peut pas être CV  
 vers  $\pi(A)$ .