

Comment nettoyer Sunstein (le bougre)

Les PDFs sont sur le github dans data/sunstein/pdf/ et les fichiers texte dans data/sunstein/processed/

Pour chaque article, partir soit du fichier texte `cermine`, soit du fichier texte `lapdf`. En général, je pars du `lapdf`, mais si vous voyez qu'il a des gros problèmes (notamment pour les articles qui ont en deuxième page une page énorme avec écrit CHICAGO et un logo en-dessous), prenez l'autre. Copiez le fichier avant de l'éditer, et la version éditée doit avoir exactement le même nom, mais avec `lapdf` ou `cermine` remplacé par `mclean`. Ouvrez l'article original en PDF à côté. Je recommande d'utiliser Notepad++ comme éditeur de texte

Note : Il y a pas besoin de lire tout le texte, notamment pour les deux derniers points de la liste, corrigez que les erreurs que vous voyez, le but c'est pas de passer des heures sur chaque article non plus.

Note 2 : Il y a des articles qui prennent très peu de temps à nettoyer, et d'autres qui peuvent prendre une heure...

Choses principales à faire :

- Vérifier chaque transition de page. Les phrases sont généralement coupées en deux à cet endroit donc il faut supprimer le saut de ligne. Attention, parfois il y a une ou plusieurs pages qui manquent du fichier texte donc il faut bien vérifier toutes les transitions une par une. Si y a des pages qui manquent, copiez-collez le texte du pdf à la place.
- Vérifier les sauts de ligne. Le but c'est qu'une ligne = un paragraphe (ou un titre de section). Il y a des endroits où l'OCR a mis une ligne de texte pour chaque ligne du pdf, donc dans ce cas il faut reconstituer les paragraphes en supprimant les sauts de ligne et en reconstituant les mots s'il y a des césures (à faire à la main! Certaines des césures sont normales si les mots contiennent naturellement un tiret, "18-year-old" est correct mais "constitu-tion" clairement pas).
- Supprimer la majorité des numéros de référence dans le texte (c'est pas grave si vous en oubliez. Ils se manifestent souvent par des chiffres à la fin des mots ou de phrases, ou des fois l'OCR les remplace des guillemets (par ex. "...in punitive damage awards.**4** This investigation...")
- Si les notes de bas de page apparaissent au milieu du texte (à la fin de chaque page), les virer
- S'il y a le temps, copier-coller les notes de bas de page à la fin du document, reconstituer les lignes correctement (une ligne par note, même s'il y a plusieurs paragraphes), et supprimer leurs numéros. Avec Notepad++ vous pouvez sélectionner en appuyant sur Ctrl pour enlever les premiers caractères de chaque ligne par exemple. Attention, il y a souvent pas mal de texte à corriger dans les notes de bas de page, vu que le texte est plus petit et l'OCR plus mauvais
- Si vous voyez des tirets entre deux mots qui servent à introduire une proposition, rajoutez un espace des deux côtés ("well-defined" → "well-defined", mais "I think-and this is true-that I like dogs." → "I think - and this is true - that I like dogs."). Si c'est des tirets longs (—) comme c'est le cas dans certains articles, rien besoin de corriger.
- S'il y a des erreurs d'OCR qui sautent aux yeux dans le texte, corriger les mots en question.
- Pour les articles avec en deuxième page CHICAGO + le logo énorme, l'OCR remplace les chiffres et les minuscules dans les titres par des caractères "boîte avec point d'interrogation ([?])" ou équivalent, c'est une erreur d'encodage, laissez-les comme ça, je les corrigerai automatiquement à la fin.

EXEMPLE 1 (TRANSITION DE PAGE)

Texte original :

21 One need not take a position on MacKinnon's broadest claims about the relationship between sexuality and sexual inequality in order to agree that the fact 1988]
22 that some women associate sexuality and violence is not a sufficient reason to permit the distribution of every film that merges sexuality and violence.

Texte nettoyé :

23 One need not take a position on MacKinnon's broadest claims about the relationship between sexuality and sexual inequality in order to agree that the fact that some women associate sexuality and violence is not a sufficient reason to permit the distribution of
24 every film that merges sexuality and violence.

EXEMPLE 2 (NOTES DE BAS DE PAGE)

Notes de bas de page copiées-collées du pdf :

31
32 1 Visiting Professor of Law, University of Chicago; Professor of Law, Osgoode Hall Law
33 School.
34 2 Professor of Law, Law School and Department of Political Science, University of Chicago.
35 The author would like to thank Mary Becker, Veronica Dougherty, Richard Fallon, Mary Ann
36 Glendon, Sara Ketchum, Larry Kramer, Michael McConnell, Frank Michelman, Geoffrey
37 Miller, Martha Minow, Richard A. Posner, Geoffrey Stone, David A. Strauss, Kathleen Sullivan,
38 and Diane Wood for valuable comments on a previous draft.
39 3 See West Coast Hotel Co. v. Parrish, 300 U.S. 379 (1937); B. ACKERMAN, RECONSTRUCT-
40 ING AMERICAN LAW (1984); L. TRIBE, AMERICAN CONSTITUTIONAL LAW 6-8 (2d ed. 1988);
41 Sunstein, Lochner's Legacy, 87 COLUM. L. REV. 873 (1987).
42 4 963 U.S. 537 (1896).
43 5 "The first unmistakably feminist voices were heard in England in the 17th century." A.
44 JAGGAR, FEMINIST POLITICS AND HUMAN NATURE 3 (1983)
45

Lignes reconstituées :

46 1 Visiting Professor of Law, University of Chicago; Professor of Law, Osgoode Hall Law School.
47 2 Professor of Law, Law School and Department of Political Science, University of Chicago. The author would like to thank Mary Becker, Veronica Dougherty, Richard Fallon, Mary Ann Glendon, Sara Ketchum, Larry Kramer, Michael McConnell, Frank Michelman, Geoffrey
48 Miller, Martha Minow, Richard A. Posner, Geoffrey Stone, David A. Strauss, Kathleen Sullivan, and Diane Wood for valuable comments on a previous draft.
49 3 See West Coast Hotel Co. v. Parrish, 300 U.S. 379 (1937); B. ACKERMAN, RECONSTRUCTING AMERICAN LAW (1984); L. TRIBE, AMERICAN CONSTITUTIONAL LAW 6-8 (2d ed. 1988); Sunstein, Lochner's Legacy, 87 COLUM. L. REV. 873 (1987).
50 4 963 U.S. 537 (1896).
51 5 "The first unmistakably feminist voices were heard in England in the 17th century." A. JAGGAR, FEMINIST POLITICS AND HUMAN NATURE 3 (1983)

Numéros supprimés :

52 Visiting Professor of Law, University of Chicago; Professor of Law, Osgoode Hall Law School.
53 Professor of Law, Law School and Department of Political Science, University of Chicago. The author would like to thank Mary Becker, Veronica Dougherty, Richard Fallon, Mary Ann Glendon, Sara Ketchum, Larry Kramer, Michael McConnell, Frank Michelman, Geoffrey
54 Miller, Martha Minow, Richard A. Posner, Geoffrey Stone, David A. Strauss, Kathleen Sullivan, and Diane Wood for valuable comments on a previous draft.
55 3 See West Coast Hotel Co. v. Parrish, 300 U.S. 379 (1937); B. ACKERMAN, RECONSTRUCTING AMERICAN LAW (1984); L. TRIBE, AMERICAN CONSTITUTIONAL LAW 6-8 (2d ed. 1988); Sunstein, Lochner's Legacy, 87 COLUM. L. REV. 873 (1987).
56 4 U.S. 537 (1896).
57 5 "The first unmistakably feminist voices were heard in England in the 17th century." A. JAGGAR, FEMINIST POLITICS AND HUMAN NATURE 3 (1983)

EXEMPLE 3 (RETOURS À LA LIGNE ET NUMÉROS DE RÉFÉRENCES)

66
67 This is a somewhat unusual subject for a lawyer; but the topic is
68 far from irrelevant to law. Some people think that the American con-
69 stitutional tradition has been punctuated by a range of constitutional
70 moments. 5 Whether or not this is so, the civil rights movement of the
71 1960s unquestionably helped transform our understandings of consti-
72 tutional principles, including most prominently rights to free speech
73 and equal protection of the laws. The political and moral claims of the
74 movement helped spur legislation' that continues to raise foundational
75 issues about our constitutional order. 7
76 Few people think that the civil rights movement actually
77 amended the constitution, 8 [...]

Texte original :

Paragraphes reconstitués :

78 This is a somewhat unusual subject for a lawyer; but the topic is far from irrelevant to law. Some people think that the American Constitutional tradition has been punctuated by a range of constitutional moments. 5 Whether or not this is so, the civil rights
79 movement of the 1960s unquestionably helped transform our understandings of constitutional principles, including most prominently rights to free speech and equal protection of the laws. The political and moral claims of the movement helped spur legislation 7
80 that continues to raise foundational issues about our constitutional order. 7
81 Few people think that the civil rights movement actually amended the constitution, 8 [...]

Numéros de références supprimés : (Note : on voit qu'il y a un 5, 7 et 8, mais pas de 6. Il faut chercher un peu, ici le 6 a été remplacé par un apostrophe après "legislation", donc il faut le supprimer aussi dans l'idéal)

EXEMPLE 4 (TRANSITION DE PAGE + NUMÉRO DE RÉFÉRENCE)

Texte original :

punitive damage awards? On the basis of
2 CHICAGO WORKING PAPER LAW ECONOMICS responses from 899 jury-eligible citizens,

Texte nettoyé :

punitive damage awards? On the basis of responses from 899 jury-eligible citizens,