

Battre son corpus tant qu'il **ait** chaud

Benoît Habert

LIR – LIMSI – CNRS & université Paris X – Nanterre

habert@limsi.fr

<http://www.limsi.fr/Individu/habert/>

## Plan

- Dégrouper automatiquement les sens
- Mesurer les changements distributionnels
- Battre son corpus : varier distance, contexte, partition
- Précautions et perspectives

## Dégrouper les sens

- Réduire/repérer polysémie et homonymie
  - une source de bruit en recherche d'information (obtention de documents non pertinents)
  - un obstacle en acquisition sémantique (« ponts » indus entre regroupements)
  - l'indice de désaccords/méconnaissances dans les débats citoyens
- Difficulté : repérer les « mots qui en cachent d'autres »
- Dimensions
  - **repérer**/caractériser
  - **corpus partitionné**/non partitionné

## Corpus et contextes

- Opposition
  - signalétique externe : LMP – 3 rubriques dans *Le Monde Parole*
  - interprétation : TCDT
    - 3 émetteurs CFTCa(vant), CFTCm(aintenue) et CFDT
    - tournants : scission entre CFDT et CFTCm ; pour la CFDT mai 68 puis 1979
- Contextes
  - syntaxiques : Syntex [Bourigault & Fabre 00]
  - « graphiques » : noms, verbes, adjectifs, adverbes dans une fenêtre
    - $\pm 5$  mots
    - phrase

## LMP ART ECO POL

<i>Rubrique</i>	<i>articles</i>	<i>mots</i>	<i>moyenne</i>	<i>maximum</i>
ART(s et médias)	2 261	1 080 620	477	2 990
ECO(nomie)	3 478	1 443 923	415	3 058
POL(itique)	2 305	1 326 576	575	5 202
Total	8 044	3 851 119		

## TCDT AFC

<i>Délimitation</i>	<i>nom</i>	<i>mots</i>
CFTC jusqu'à la scission de 1964 incluse et CFDT d'avant mai 1968	TC45-64_DT65-67	29 704
CFTC maintenue	TC65-90	31 673
CFDT « radicalisée »	DT70-76	25 596
CFDT « recentrée »	DT79-92	34 677
Total (CFTC CFDT 1945-1992)	TCDT	121 650

## Relations produites par Syntex 1/2

TCDT Phrase de la résolution générale CFDT 1965

Depuis la mise en place, en 1950, de ce **salaire minimum interprofessionnel garanti en dessous** duquel aucun **travailleur** de plus de 18 ans ne doit **être payé**, la **protection sociale** qui **présidait à cette institution s'est sérieusement dégradée.**

## Relations produites par Syntex 2/2

<i>Lemme 1</i>	<i>POS 1</i>	<i>relation</i>	<i>lemme 2</i>	<i>POS 2</i>
salaire	N	EPI	minimum	A
salaire	N	EPI	interprofessionnel	A
salaire	N	EPI	garantir	A
garantir	V	en	dessous	N
payer	V	OBJ	travailleur	N
se dégrader	V	SUJ	protection	N
se dégrader	V	SUJ	protection social	S
présider	V	SUJ	protection	N
présider	V	SUJ	protection social	S
protection	N	EPI	social	A
présider	V	à	institution	N



## Proximités distributionnelles

V dans <i>OBJ__travailleur</i>	partie	#	N dans <i>OBJ_payer</i>	partie	#
concerner	DT70-76	1	travailleur	DT70-76	1
défendre	DT70-76	1	travailleur	TC45- 64_DT65-67	2
garantir	TC45- 64_DT65-67	1	action	TC45- 64_DT65-67	1
maintenir	DT70-76	1			
payer	DT70-76	1			
payer	TC45- 64_DT65-67	2			
rendre	TC45- 64_DT65-67	1			

## « Colorer » les occurrences selon la partie

- Une occurrence dans une partie est suffixée du nom de cette partie  
*travailleur* dans TCDT AFC  $\longrightarrow$  *travailleur\_\_TC45-64\_DT65-67*,  
*travailleur\_\_DT70-76*, *travailleur\_\_DT79-92* et *travailleur\_\_TC65-90*
- On peut mélanger sans les confondre ces hétérographes artificiels et leurs contextes

### Exemple

<i>Hétérographe artificiel</i>	<i>contexte</i>	<i>fréquence</i>
<i>travailleur__TC45-64_DT65-67</i>	<i>OBJ__payer</i>	1
<i>payer__TC45-64_DT65-67</i>	<i>OBJ__travailleur</i>	1

## Visualiser/mesurer les changements de distribution

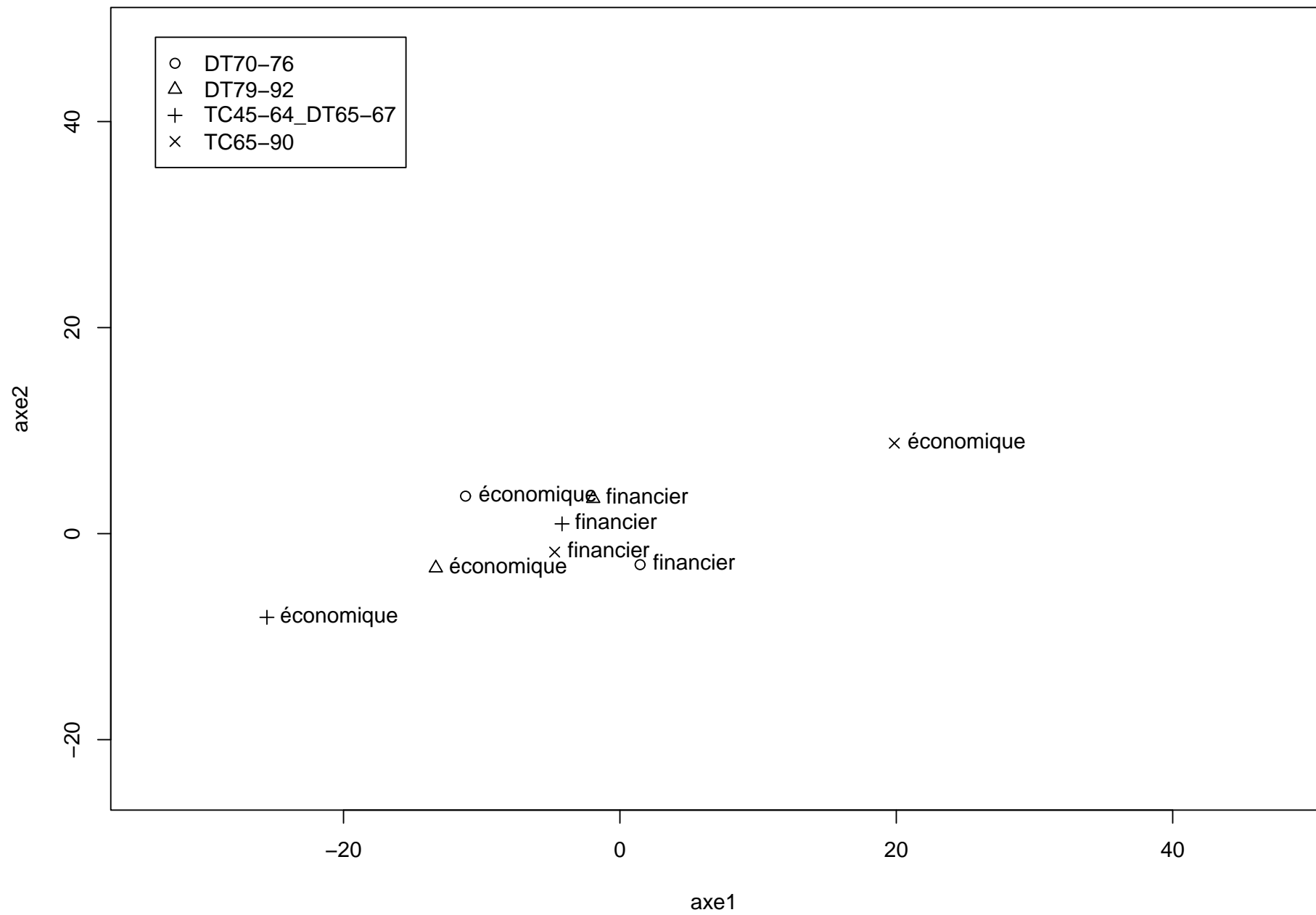
- Visualisation : projection des hétérographes dans l'espace des contextes
    - proximité *financier* ( $\approx$  ? convergence)
    - dispersion *économique* ( $\approx$  ? divergence)
    - « bande à part » *travailleur, réel*
  - Distance cumulée entre hétérographes d'un même lemme
- 2 distances
- cosinus + pondération des traits (tf.idf)
  - Jaccard (sans pondération)

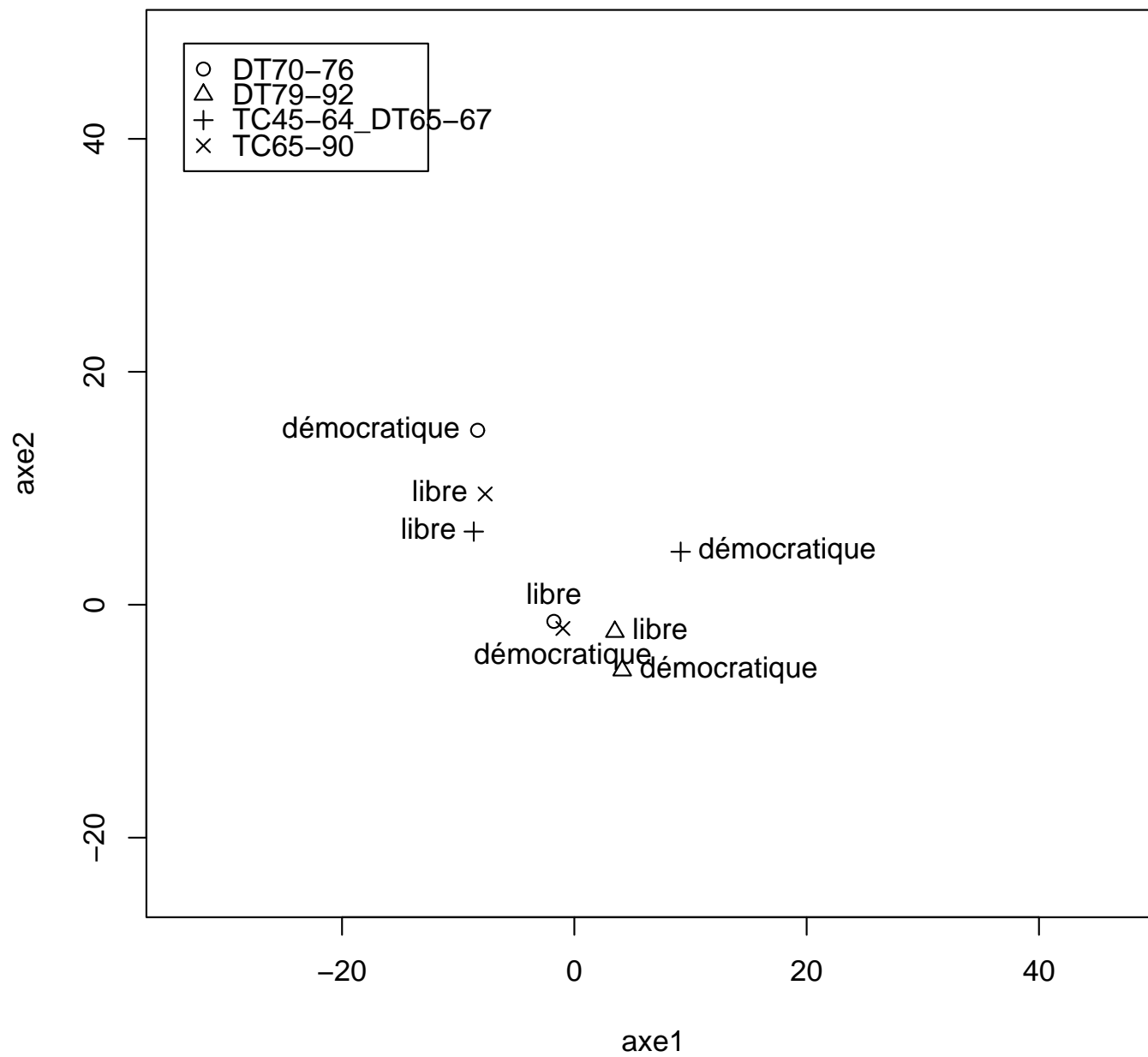
## TCDT distances (cosinus, Jaccard) des hétérographes de *travailleur*

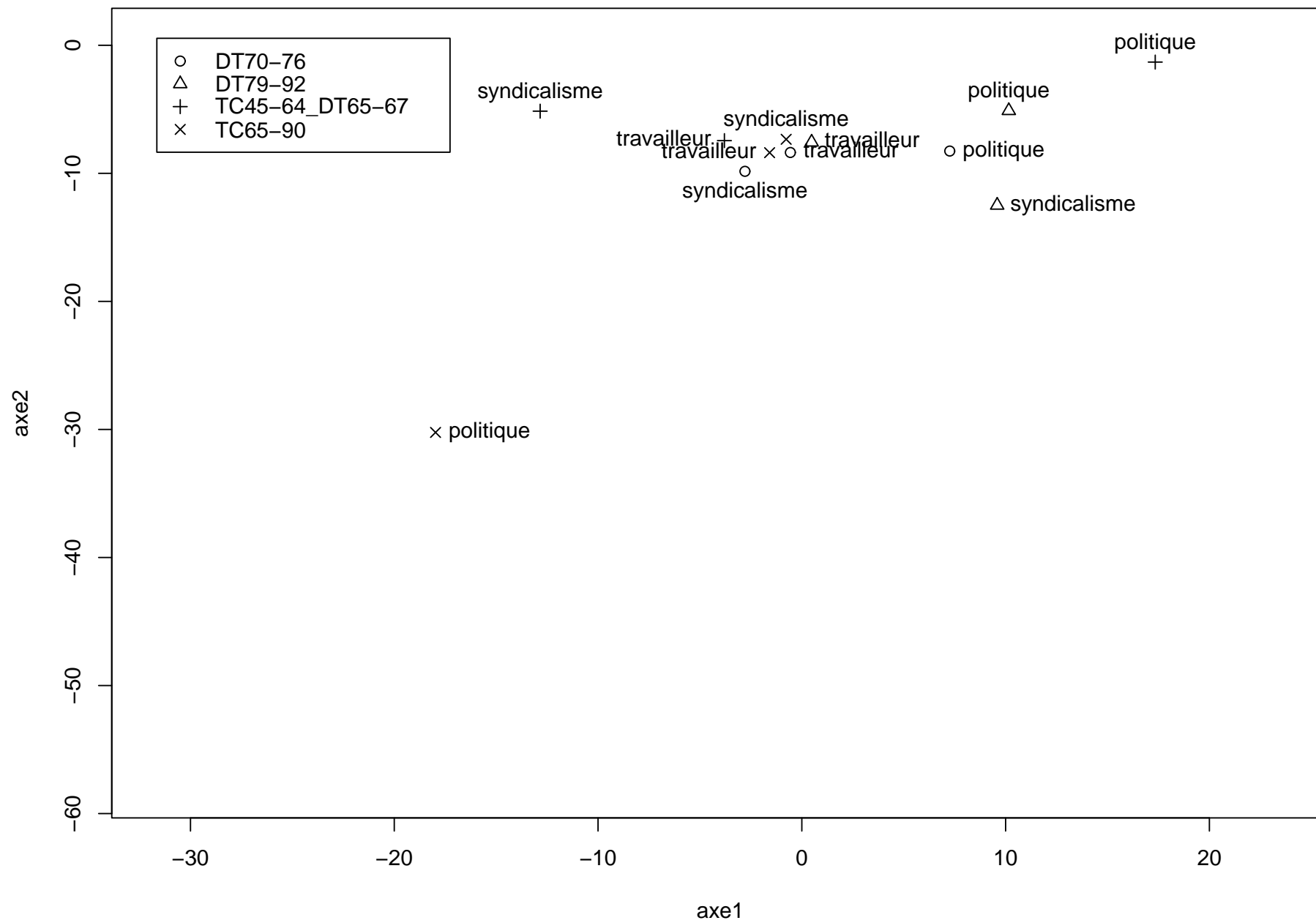
<i>Hétérographe 1</i>	<i>hétérographe 2</i>	<i>cosinus</i>	<i>Jaccard</i>
~__DT70-76	~__DT79-92	0.93	0.949
~__DT70-76	~__TC45-64_DT65-67	0.902	0.958
~__DT70-76	~__TC65-90	0.936	0.935
~__DT79-92	~__TC45-64_DT65-67	0.972	0.986
~__DT79-92	~__TC65-90	0.973	0.957
~__TC45-64_DT65-67	~__TC65-90	0.83	0.93
Somme des distances pour <i>travailleur</i>		5.543	5.715

## TCDT 10 premiers N partagés par distance cumulée entre hétérographes

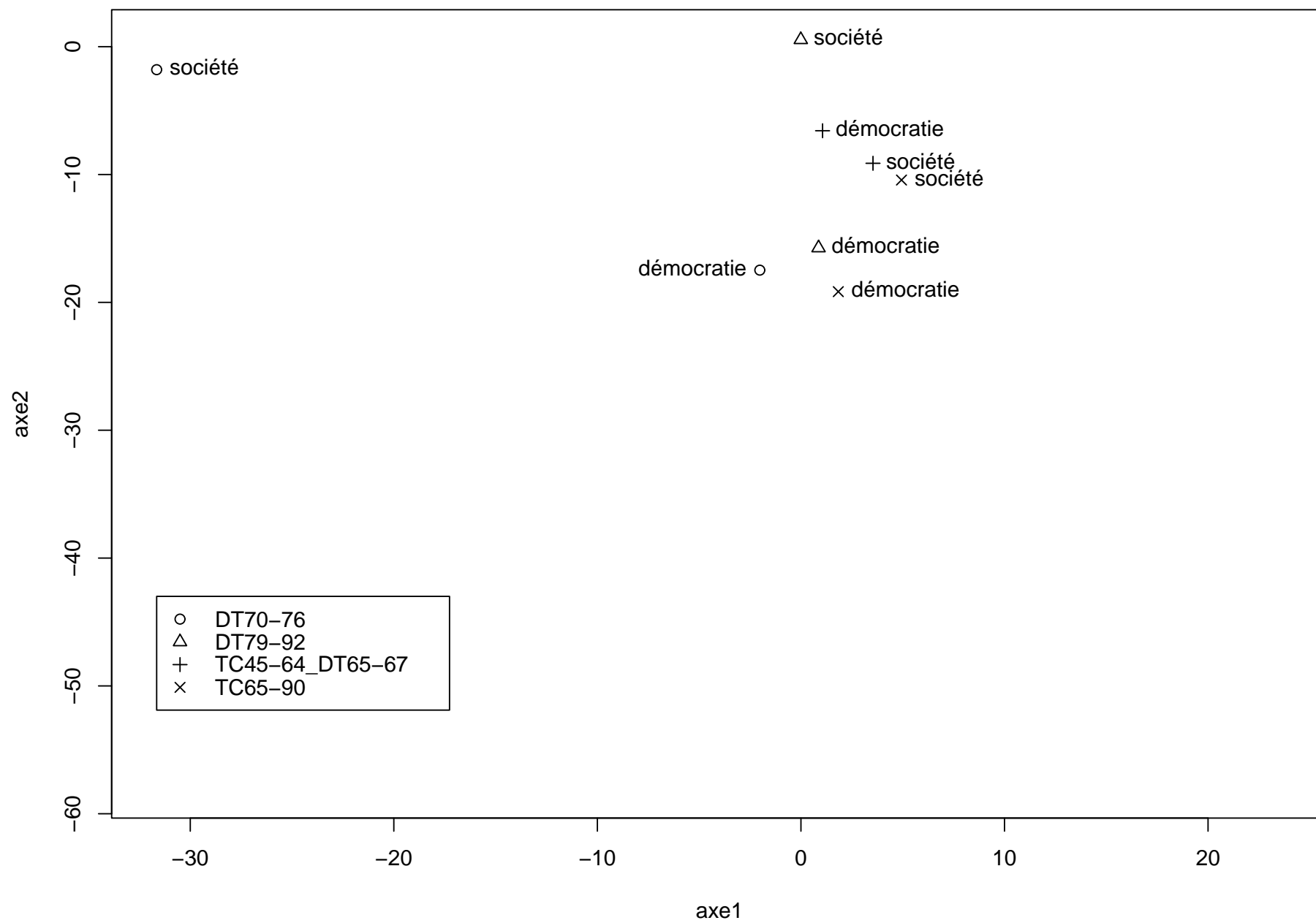
<i>Lemme</i>	<i>cosinus</i>	<i>Jaccard</i>
ensemble	5.659	5.695
syndicalisme	5.623	5.734
travailleur	5.543	5.715
institution	5.515	5.68
relation	5.391	5.311
intérêt	5.34	5.564
solidarité	5.335	5.327
condition	5.329	5.432
situation	5.326	5.516
niveau	5.256	5.696











## Résultats LMP (Syntex + cosinus)

<i>POS</i>	<i>lemmes</i>
A	même, célèbre, classique, indispensable, fameux, possible, issu, traditionnel, présent, directeur, japonais, italien, militaire, complet, indépendant, urbain, allemand, populaire, permanent, rare
N	résultat, retour, titre, film, voix, salle, construction, auteur, évolution, autorité, cours, ouverture, gestion, demande, image, mouvement, maison, dirigeant, vente, ligne
V	remplacer, accompagner, découvrir, regrouper, naître, demeurer, apparaître, autoriser, continuer, se retrouver, entraîner, sembler, risquer, envisager, lier, construire, considérer, travailler, conduire, diriger

## Résultats TCDT AFC (Syntex + cosinus)

<i>POS</i>	<i>lemmes</i>
A	véritable, nécessaire, financier, actuel, européen, essentiel, international, libre, politique, confédéral, mondial, tout, national, démocratique, nationale, économique, réel, local, professionnel, syndicaux
N	ensemble, syndicalisme, travailleur, institution, relation, intérêt, solidarité, condition, situation, niveau, moyen, société, effort, démocratie, entreprise, responsabilité, objectif, activité, emploi, politique
V	assurer, renforcer, devoir, confirmer, contribuer, développer, ouvrir, améliorer, assumer, réduire, apporter, poser, répondre, poursuivre, participer, réaffirmer, affirmer, continuer, engager, réaliser

## Variation distributionnelle $\neq$ divergence

- Variation des contextes de verbes et d'adjectifs, voire de noms, ne correspondant pas à des différences d'acceptions
- Adjectifs
  - « de rangement » *social, économique, européen, mondial*
  - de nationalité *japonais, italien, allemand*
  - évaluatifs *indispensable, fameux, célèbre, véritable*
- Verbes
  - copules *sembler, apparaître, demeurer*
  - de parole *mandater, réaffirmer, rappeler*
  - dénotant le changement *poursuivre, améliorer, développer*
- Noms « outils » *niveau, ensemble, situation, moyen*

## Variations dans LMP 1/2

- Artefact : distance cumulée forte par emploi écrasant dans une partie *film* (ART : 1 200 contextes / 25 POL / 21 ECO)
  - ART *auteur, salle*
  - ECO *cours*
  - POL *mouvement*
- Emploi prépondérant (mais déséquilibré) dans 2 parties
  - ECO > POL *autorité, construction, demande, évolution, gestion, résultat*
  - POL > ECO *dirigeant*
  - ART > POL *image, voix*
  - ART > ECO *maison*
  - ECO > ART *ouverture, titre, vente*

## Variations dans LMP 2/2

- Spécialisations
  - *cours* monnaie, Bourse (ECO) / processus (ART)
  - *mouvement* collectif (POL) / déplacement (ART) et évolution (ART, ECO)
  - *image* concret / abstrait (ART donne plus de place au premier sens que POL et ECO)
  - *ouverture* début d'un morceau de musique (ART)
  - *voix* qualité, tessiture (ART) / votes (POL)

## Variations dans TCDT AFC 1/3

- Unité en thème et en genre du corpus → distributions moins iniques que LMP
- Associations dominantes
  - DT79-92 et TC65-90 *activité, emploi, entreprise*
  - TC45-64\_DT65-67 et DT70-76 *démocratie*
  - DT70-76 et DT79-92 *société*
- *condition*
  - contextes peu révélateurs *EPI\_\_fondamental, OBJ\_\_améliorer, EPI\_\_favorable, EPI\_\_indispensable, EPI\_\_essentiel, EPI\_\_nécessaire*
  - opposition entre DT70-76 - DT79-92 OBJ (c'est quelque chose sur laquelle agir) et TC65-90 - TC45-64\_DT65-67 avant tout EPI ou SUJ

## Variations dans TCDT AFC 2/3

- *démocratie*

TC45-64\_DT65-67 / DT70-76 renvoi à une menace *OBJ\_restaurer*,  
*OBJ\_rétablir*

- *société*

- partages à 4 *EPI\_\_démocratique* ou 3 *OBJ\_\_construire*
- TC45-64\_DT65-67 et TC65-90 *pour\_\_homme*
- DT70-76 et DT79-92 *EPI\_\_autogestionnaire*, *EPI\_\_capitaliste*,  
*EPI\_\_inégalitaire*, *EPI\_\_socialiste* et *OBJ\_\_changer*
- DT70-76 contextes très politiques *EPI\_\_aliénant*, *EPI\_\_autogérer*,  
*EPI\_\_autogéré*, *EPI\_\_décentraliser*, *EPI\_\_futur*,  
*EPI\_\_hiérarchique*, *EPI\_\_technocratique*, à base de *\_\_autogestion*
- DT79-92 *EPI\_\_ouvert*, *EPI\_\_pluraliste*, *EPI\_\_pluriculturel*,  
*EPI\_\_répressif*, *EPI\_\_éclaté*



## Variations dans TCDT AFC 3/3

- *syndicalisme*
  - sauf TC65-90 *EPI\_\_démocratique*
  - TC65-90 et TC45-64\_DT65-67 *EPI\_\_chrétien, EPI\_\_libre*
  - DT70-76 peu de contextes propres
  - TC65-90 - TC45-64\_DT65-67 SUJ privilégié par rapport à DT79-92.
  - échos *EPI\_\_constructif* TC45-64\_DT65-67 / *EPI\_\_responsable* TC65-90

- *travailleur*
  - *SUJ\_participer* 3 parties sur 4
  - TC45-64\_DT65-67 et TC65-90 *EPI\_chrétien* et *EPI\_manuel*
  - TC65-90 faible emploi net
  - DT79-92 fragmentation de la catégorie *EPI\_français*,  
*EPI\_immigré*, *EPI\_originaire*, *EPI\_syndiqué*, *EPI\_permanent*,  
*EPI\_engagé*
  - DT70-76 représentation plus conflictuelle *OBJ\_défendre*, *en\_lutte*,  
*en\_situation\_de\_subordination*, *face\_à\_pouvoir*.
  - DT79-92 et DT70-76 favorisent SUJ sur OBJ.

## Varier les distances – LMP

<i>Distance</i>	
cosinus	résultat, retour, titre, <b>film</b> , voix, salle, construction, auteur, évolution, autorité, cours, ouverture, gestion, demande, image, mouvement, maison, dirigeant, vente, ligne
Jaccard	musique, <b>film</b> , théâtre, voix, musée, hausse, droite, femme, reprise, compagnie, oeuvre, achat, chaine, scène, histoire, cours, art, actionnaire, opposition, maire

## Varier les distances – TCDT AFC

<i>Distance</i>	
cosinus	<b>ensemble, syndicalisme, travailleur, institution, relation, intérêt, solidarité, condition, situation, niveau, moyen, société, effort, démocratie, entreprise, responsabilité, objectif, activité, emploi, politique</b>
Jaccard	<b>syndicalisme, travailleur, emploi, niveau, place, ensemble, institution, service, objectif, moyen, société, pouvoir, responsabilité, entreprise, effort, intérêt, mise, démocratie, droit, prise</b>

## Varier les contextes

- Version du corpus étiquetée par/pour Syntex
- Restriction aux noms, adjectifs, verbes et adverbes (degré, nég.)
- Fenêtre graphique (ie sans relations syntaxiques)
- 2 versions
  - $\pm 5$  mots (quelle que soit leur catégorie)
  - phrase
- Exemple pour *travailleur* dans la phrase-exemple
  - Syntex {*payer*/OBJ}
  - $\pm 5$  mots {**dessous**}
  - phrase {mise, en\_place, 1950, salaire, minimum, interprofessionnel, garantir, **dessous**, an, *payer*, protection, sociale, présider, institution, se\_dégrader, sérieusement}

Les matrices se remplument

	Syntax	$\pm 5$ N, A, V, Adv	Phrase N, A, V, Adv
% occupé	0.25	0.61	2.94
Formes	271	438	438
Hétérographes	1 084	1 748	1 748
Fmax	232	692	4 217
moyenne	12.58	35.47	226
1er quart.	3.0	7.0	42.0
2e quart.	7.0	18.0	113.0
3e quart.	15.0	41.0	260.25
Contextes	3 648	4 318	4 746
Fmax	405	1163	6 517
moyenne	3.73	14.36	83.33

## TCDT 4 parties – Contextes partagés par toutes les parties pour *travailleur*

Syntex – 0	
$\pm$ 5 N, A, V, Adv	
<b>ne ; pays ; pouvoir ; être</b>	4
Phrase N, A, V, Adv	
<p>action ; activité ; apporter ; assurer ; autre ; besoin ; collectif ; condition ; congrès ; croissance ; droit ; développement ; emploi ; en particulier ; entreprise ; faire ; grand ; immigré ; intérêt ; <b>ne</b> ; niveau ; nouveau ; nécessaire ; nécessité ; organisation ; participer ; pas ; <b>pays</b> ; plus ; politique ; <b>pouvoir</b> ; public ; responsabilité ; rôle ; salaire ; salarié ; secteur ; situation ; social ; structure ; syndicalisme ; tout ; travail ; travailleur ; vie ; volonté ; économie ; économique ; <b>être</b></p>	49



TCDT 4 parties – Nombre de contextes propres pour *travailleur*

	Syntex	$\pm 5$	Phrase
DT70-76	24	67	159
DT79-92	26	59	147
TC45-64_DT65-67	36	94	374
TC65-90	11	41	145

TCDT 4 parties – Syntex – Nombre de contextes partagés pour *travailleur*

DT70-76	DT79-92	TC45-64_DT65-67	TC65-90
<i>24</i>	<i>26</i>	<i>36</i>	<i>11</i>
2	2		
2		2	
1			1
	1	1	
	1		1
		3	3
	0	0	0
1		1	1
1	1		1
0	0	0	
0	0	0	0

TCDT 4 parties –  $\pm 5$  N, A, V, Adv – Partages *travailleur*

DT70-76	DT79-92	TC45-64_DT65-67	TC65-90
67	59	94	41
15	15		
7		7	
5			5
	8	8	
	2		2
		9	9
	4	4	4
3		3	3
3	3		3
7	7	7	
4	4	4	4

TCDT 4 parties – Phrase N, A, V, Adv – Partages *travailleur*

DT70-76	DT79-92	TC45-64_DT65-67	TC65-90
<i>159</i>	<i>147</i>	<i>374</i>	<i>145</i>
<i>53</i>	<i>53</i>		
<i>54</i>		<i>54</i>	
<i>18</i>			<i>18</i>
	<i>56</i>	<i>56</i>	
	<i>12</i>		<i>12</i>
		<i>71</i>	<i>71</i>
	<i>27</i>	<i>27</i>	<i>27</i>
<i>27</i>		<i>27</i>	<i>27</i>
<i>9</i>	<i>9</i>		<i>9</i>
<i>43</i>	<i>43</i>	<i>43</i>	
<i>49</i>	<i>49</i>	<i>49</i>	<i>49</i>

Comparaison des distances et par type de contexte



<i>Contextes</i>	<i>lemmes des k premiers N par distance cumulée (cosinus) entre hétérographes</i>
Syntax	ensemble, syndicalisme, travailleur, institution, relation, intérêt, solidarité, condition, situation, niveau, moyen, société, effort, démocratie, entreprise, responsabilité, objectif, activité, emploi, politique
± 5	<i>fait, intervention, exigence, en_vue, occasion, domaine, base, accroissement, propriété, projet, perspective, insertion, institution, victime, investissement, réalisation, renforcement, système, état, utilité</i>
Phrase	<b>histoire, assistance, occasion, intégration, dépendance, appareil, aspect, pensée, dégradation, sauvegarde, esprit, liaison, jour, utilité, conjoncture, libéralisme, intéressé, connaissance, adhésion, nature</b>

Contextes	<i>lemmes des k premiers N par distance cumulée (Jaccard) entre hétérographes</i>
Syntax	syndicalisme, travailleur, emploi, niveau, place, ensemble, institution, service, objectif, moyen, société, pouvoir, responsabilité, entreprise, effort, intérêt, mise, démocratie, droit, prise
± 5	revenu, <i>intervention</i> , <i>exigence</i> , perspective, chômage, part, <i>fait</i> , <i>réalisation</i> , réduction, <i>projet</i> , création, enseignement, propriété, <i>investissement</i> , institution, accroissement, sécurité, <i>en_vue</i> , nécessité, régime
Phrase	insuffisance, <b>histoire</b> , égard, capital, règle, licenciement, esprit, poursuite, <b>appareil</b> , disposition, <b>liaison</b> , <b>jour</b> , <b>dépendance</b> , <b>pensée</b> , nationalisation, <b>aspect</b> , communauté, <b>intégration</b> , <b>adhésion</b> , <b>assistance</b>

## Varier les partitions

- Contrastes d'émetteurs : CFTCa(vant), CFTCm(aintenue), CFDT
  - CFTCa / CFTCm / CFDT
  - CFTCa / CFTCm
  - CFTCa / CFDT
  - CFTCm / CFDT
- En aval d'une interprétation : 4 parties issues de l'interprétation d'une analyse factorielle des correspondances

Résultats pour d'autres partitions

	<p><i>Partition lemmes des k premiers N par distance (cosinus) entre hétérographes</i>  (a : avant scission, m : maintenue – contexte : <math>\pm 5</math> N, A, V, Adv)</p>
TCa DT	perspective, partie, libération, exploitation, contact, compensation, productivité, collectivité, patronat, instauration, valeur, discussion, projet, effet, demande, profession, nature, adhésion, enfant, an
TCa TCm DT	partie, direction, nature, libération, combat, discussion, rémunération, projet, intervention, réalité, instauration, système, régression, statut, difficulté, autonomie, valeur, revalorisation, soutien, en_vue
TCa TCm	statut, généralisation, forme, an, partie, système, accroissement, réglementation, baisse, souci, direction, commission, syndicat, intervention, information, en_vue, nationalisation, contexte, méthode, devoir
TCm DT	obligation, lien, dialogue, coordination, présence, modalité, discussion, conception, source, partie, environnement, risque, égard, nationalisation, direction, argent, nature, concertation, thème, manière

	<i>Partition lemmes des k premiers N par distance (Jaccard) entre hétérographes</i> (a : avant scission, m : maintenue – contexte : $\pm 5$ N, A, V, Adv)
TCa DT	pratique, perspective, adhérent, projet, patronat, capacité, part, temps, accord, autonomie, réalité, mobilisation, fonctionnement, exploitation, égalité, libération, profit, difficulté, collectivité, choix
TCa TCm DT	rémunération, partie, marché, direction, transformation, tendance, valeur, autonomie, accord, demande, temps, projet, répartition, intervention, grève, forme, combat, union, productivité, an
TCa TCm	accord, statut, syndicat, an, initiative, caractère, information, crise, rémunération, en_vue, marché, loi, temps, équipement, négociation, généralisation, forme, employeur, dépense, esprit
TCm DT	classe, conception, unité, proposition, fonctionnement, union, rémunération, mouvement, cotisation, coordination, discussion, application, fonds, débat, discrimination, conscience, adhérent, valeur, risque, reconnaissance

## Caractériser les distributions

- Voisins dans la partie
  - changement des traits (ensemble des contextes de la partie pour la catégorie / ensemble des contextes partagés par les lemmes retenus)
  - les voisinages peuvent être divergents
- Examen exploratoire (AFC) des ensembles de contextes pour repérer des directions ( $\neq$  acceptions)
- Examen exploratoire (AFC) du regroupement des ensembles de contextes selon les parties
- Partitionnement en  $k$  (arbitraire) classes des ensembles de contextes, puis examen (AFC) des regroupements
  - « individus »
    - les ensembles de contextes (ici la phrase)
    - les classes d'ensembles

Caractériser par les voisins : TCDT AFC *travailleur*



<i>Partie</i>	<i>Jaccard</i>
TC45-64_DT65-67	syndicat (0.039), organisation (0.031), <b>diplômé</b> (0.024), <b>étape</b> (0.024), <b>civilisation</b> (0.024), <b>ouvrier</b> (0.024), potentiel (0.023), <b>voie</b> (0.023), <b>morale</b> (0.023), concurrence (0.023).
DT70-76	ensemble (0.045), <b>prédétermination</b> (0.032), <b>prédominance</b> (0.032), <b>emprise</b> (0.032), <b>participe</b> (0.031), <b>cours</b> (0.031), adhésion (0.031), hégémonie (0.031), <b>militaire</b> (0.031), <b>cohérence</b> (0.031).
DT79-92	<b>économie</b> (0.062), mouvement (0.048), salarié (0.036), <b>imagination</b> (0.033), degré (0.032), <b>diversification</b> (0.032), refus (0.032), <b>masse</b> (0.032), particulier (0.031), réalisation (0.031).
TC65-90	jeu (0.083), syndicalisme (0.07), service (0.066), travail (0.057), entreprise (0.057), <b>prévention</b> (0.056), <b>isme</b> (0.056), <b>père</b> (0.056), établissement (0.053), <b>morale</b> (0.053).

Voisins et distances : TCDT AFC *travailleur*

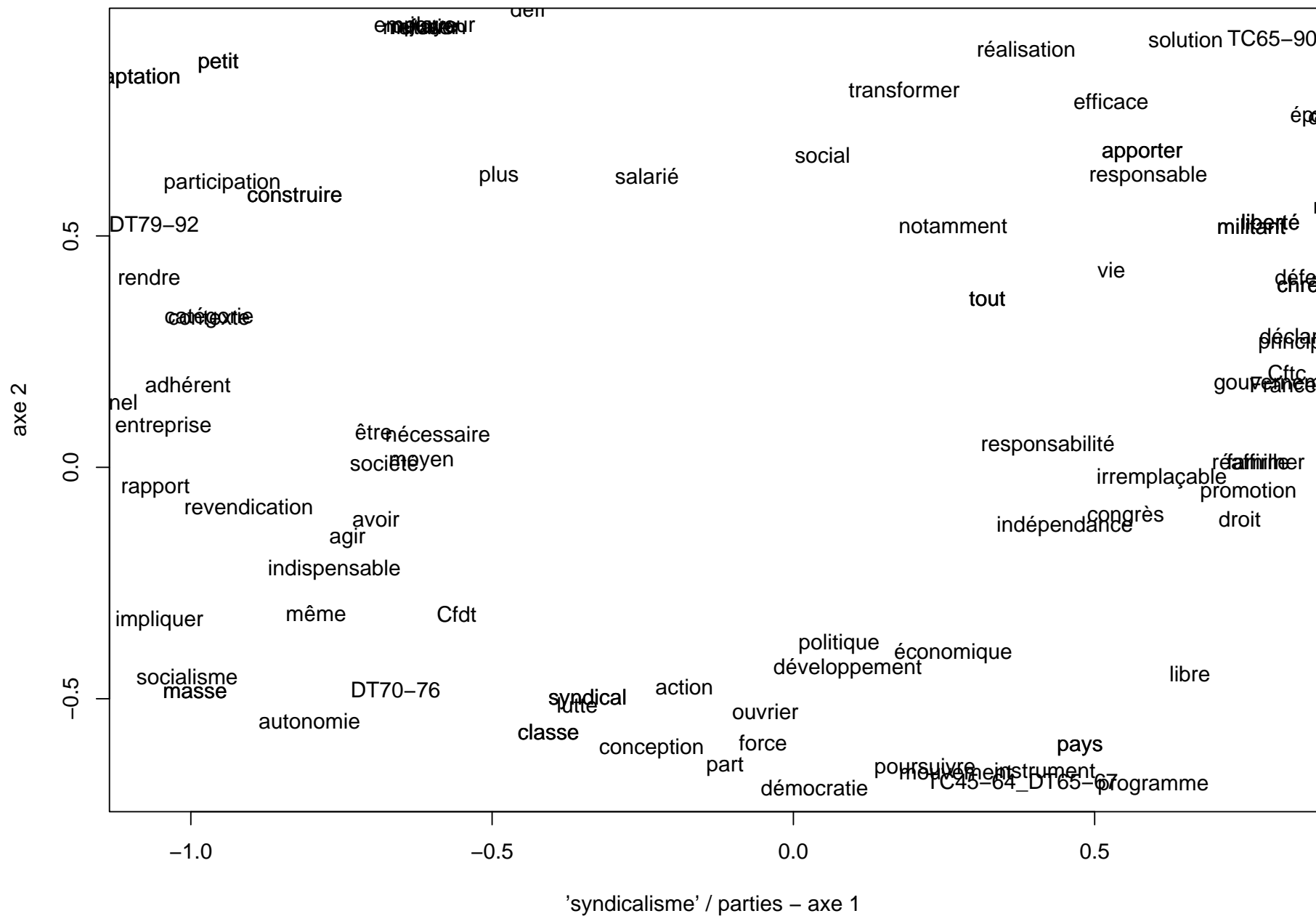
<i>Partie</i>	<i>Cosinus</i>	<i>Jaccard</i>
TC45- 64_DT65- 67	<b>civilisation</b> (0.187), <b>étape</b> (0.153), syndicalisme (0.139), personne (0.137), <b>morale</b> (0.133), <b>concurrency</b> (0.115), <b>diplômé</b> (0.113), <b>voie</b> (0.109), calcul (0.108), <b>ouvrier</b> (0.101).	syndicat (0.039), organisation (0.031), <b>diplômé</b> (0.024), <b>étape</b> (0.024), <b>civilisation</b> (0.024), <b>ouvrier</b> (0.024), potentiel (0.023), <b>voie</b> (0.023), <b>morale</b> (0.023), <b>concurrency</b> (0.023).

## Explorations autour de *travailleur*

- AFC des 4 427 associations (1 408 contextes différents) de *travailleur* avec les noms, adjectifs, verbes et adverbes des 252 phrases où il figure (occupation : 1.19%, moy. phrase : 19.50, moy. contexte : 3.14)
- AFC des 4 parties (occupation : 38.37%, moy. partie : 1 106.75, moy. contexte : 3.14)
  - TC45-64\_DT65-67 77 phrases, 1 929 associations (817 contextes)
  - DT70-76 67 phrases, 894 associations (462 contextes)
  - DT79-92 67 phrases, 959 associations (474 contextes)
  - TC65-90 41 phrases, 648 associations (408 contextes)
- AFC des  $k$  classes fournies par un classifieur

*travailleur* par parties

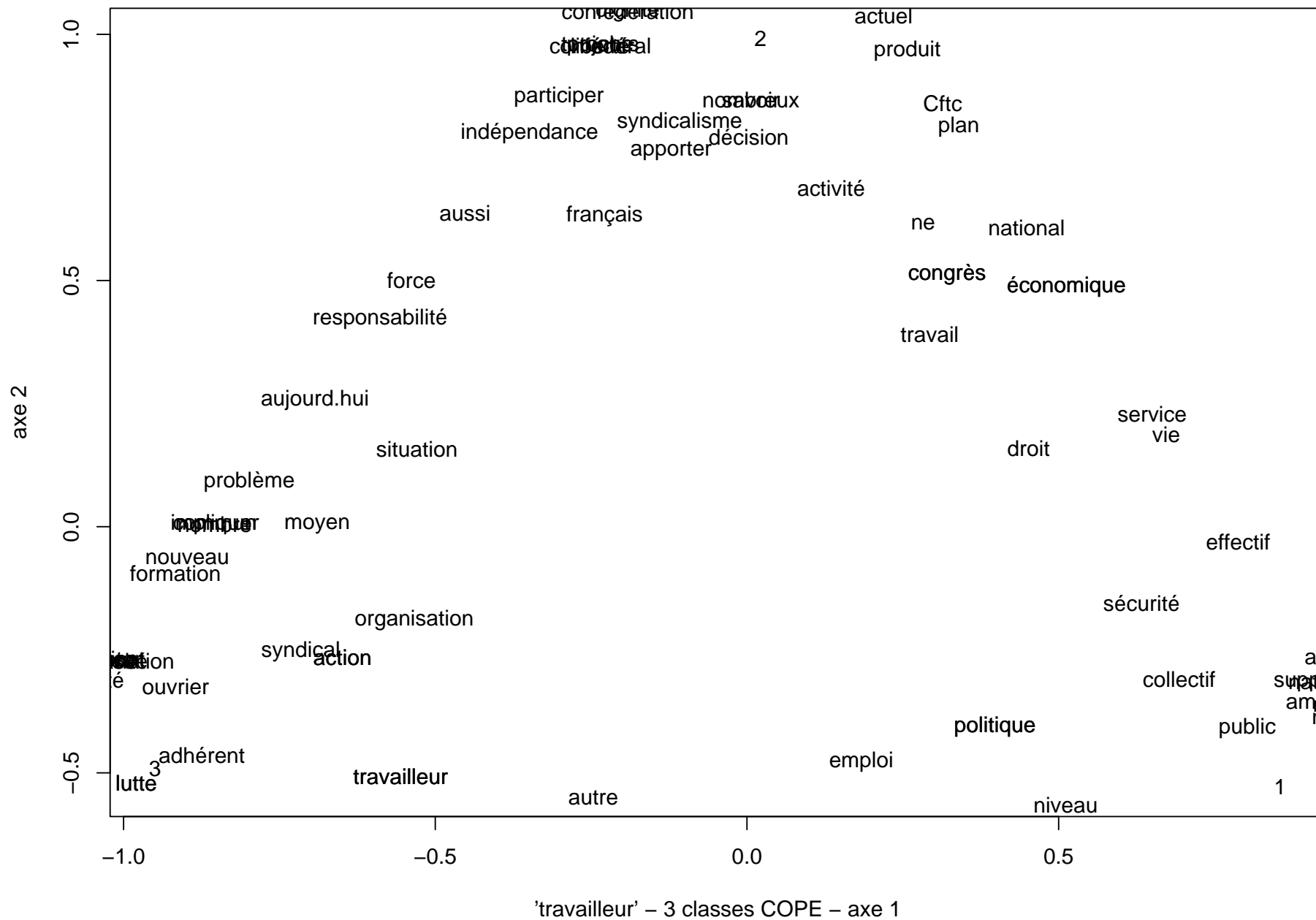
### TCDD 4 parties – AFC N V A R dans phrase



3 classes pour *travailleur*

contraste entre 3 classes pour *travailleur*

TCDT – AFC N V A R dans phrase





## Se distancier des distances

- Pondérer en fonction de la taille du corpus et de sa redondance (LMP/TCDT : faible taille + évolution idéologique sur 47 ans)
- Varier les distances pour contrebalancer la compulsion interprétative
- Distance cumulée des distributions  $\neq$  différence d'acceptation

## Pallier l'émiettement des contextes

- Matrices « fantomatiques »
  - Fragmentation supplémentaire liée aux hétérographes artificiels
  - —→ points de contact limités entre hétérographes artificiels
- Multiplier et « alourdir » les partages de contextes
  - Associer contextes syntaxiques/graphiques
  - Recherche ciblée de contextes (extraction d'information)
- Réduction de l'espace des traits
  - Calcul de similarités de second ordre
  - Projection de connaissances exogènes (synonymes, classes de mots)

## Continuer à battre ses corpus

- Amélioration par rapport à une visualisation de la dispersion des hétérographes artificiels
- Objectif du repérage de « mots qui en cache d'autres » non atteint
  - instabilité des distances → ordre fourni peu fiable
  - « brouillage » des décalages trop importants des répartitions dans les parties du corpus
  - plutôt inflexions de sens qu'écart d'acceptions
- Ne pas attendre l'obésité des corpus : les corpus à « divergences » (forums, débats citoyens, controverses) ne sont pas forcément volumineux