



APPORTS DU *WEB-SCRAPING* POUR UN SUIVI À HAUTE FRÉQUENCE DU MARCHÉ IMMOBILIER : LE CAS DU ROYAUME-UNI

CNIS – COMMISSION SYSTÈME FINANCIER ET FINANCEMENT DE L'ÉCONOMIE

27 MAI 2021

J.-C. BRICONGNE
BANQUE DE FRANCE

Messages principaux :

- Les données issues du *web-scraping* (téléchargement massif de données sur Internet) permettent un suivi granulaire [tout en veillant à la représentativité] et plus régulier des marchés immobiliers, et offrent un angle complémentaire, avec le point de vue des vendeurs
- Dans le cas du Royaume-Uni, le marché a été gelé pendant le premier confinement, avec un effondrement temporaire des offres postées, et une relative stabilité des prix (absolus) postés sauf sur Londres
- Les données de *scraping* sont cohérentes avec les sources officielles et permettent d'obtenir des prix en niveau, qui apportent de l'information par rapport aux indices, et permettent de mieux prévoir les retournements et d'estimer les stocks de richesse immobilière





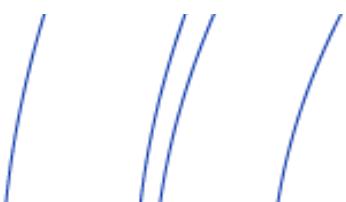
Nombreuses informations publiques en temps réel sur les sites immobiliers au UK :
92 % des agences immobilières publient des annonces sur Internet

Exemples d'utilisation de données alternatives pour suivre le marché immobilier :
Kulkarni *et al.*, 2009 avec les *Google Trends* ou *web-scraping* (par exemple Hanson et Santas, 2014 ou Bricongne *et al.*, 2019)

Le *web-scraping* offre le point de vue des vendeurs (par opposition avec les transactions finales qui résultent de l'interaction entre vendeurs et acheteurs)

Couverture du marché UK dans cette période particulière de Covid-19 & Brexit, sachant que l'approche peut être étendue à d'autres zones géographiques

Couverture des sites utilisés variable selon les régions, mais assez élevée dans tous les cas pour assurer une bonne représentativité au niveau infra-national





Données récupérées par *web-scraping* des sites majeurs d'annonces immobilières au UK : Rightmove, Zoopla et OnTheMarket.

Pour améliorer la couverture sur des régions plus spécifiques, PropertyPal – un site immobilier spécialisé dans l'Irlande du Nord – et S1Homes, son *alter ego* pour l'Écosse, sont également couverts.

Grâce à ces sites *Web*, plus de 1,5 million d'offres immobilières sont téléchargées en moyenne chaque jour. Environ les deux tiers de ces dernières sont des offres de vente (cf. tableau). Les données pour Zoopla sont téléchargées (et nettoyées) depuis début mars 2020 alors que les données d'autres sites *Web* le sont depuis juillet 2020.

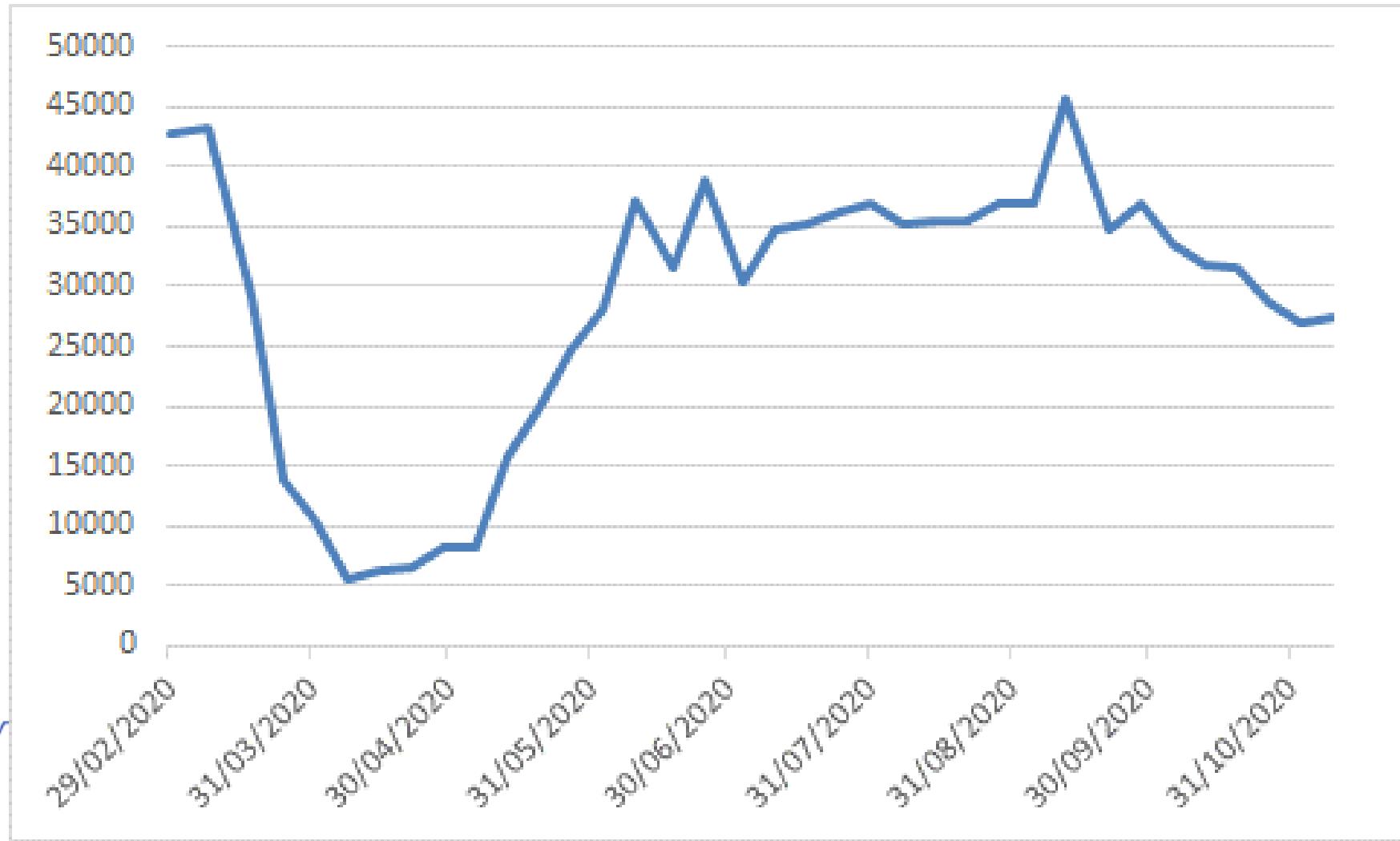




Nombre moyen d'offres *web-scrapées* par jour (par URL unique)

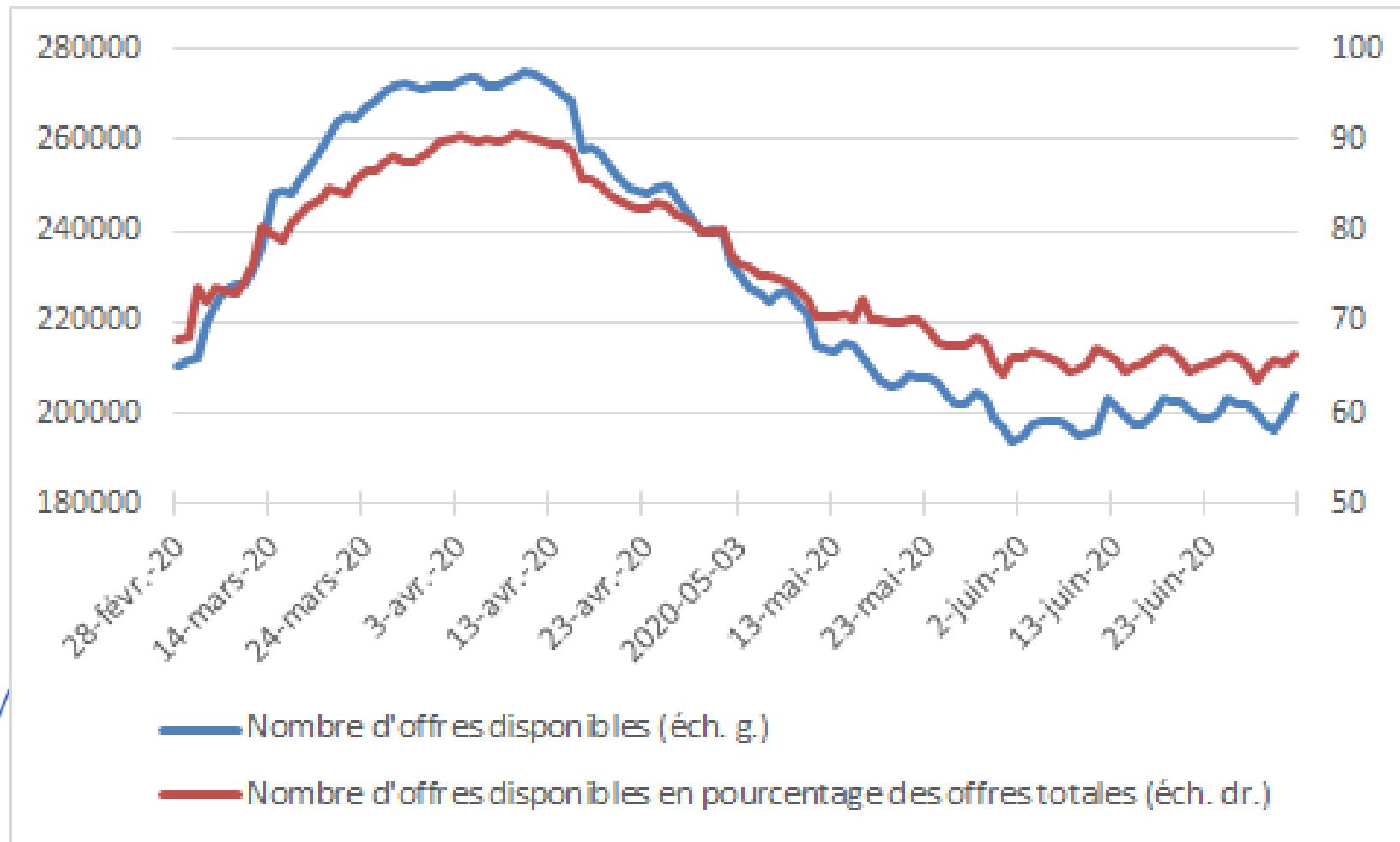
Nom du site	Vente (résidentiel)	Location (résidentiel)	Vente (commercial)	Location (commercial)
Zoopla	300 000	200 000	15 000	35 000
Rightmove	500 000	200 000	20 000	25 000
OnTheMarket	250 000	100 000	5 000	10 000

Des nouvelles offres immobilières hebdomadaires en baisse très forte lors du premier confinement et moindre lors du deuxième (en unités/semaine)



Sources : Zoopla calculs des auteurs.

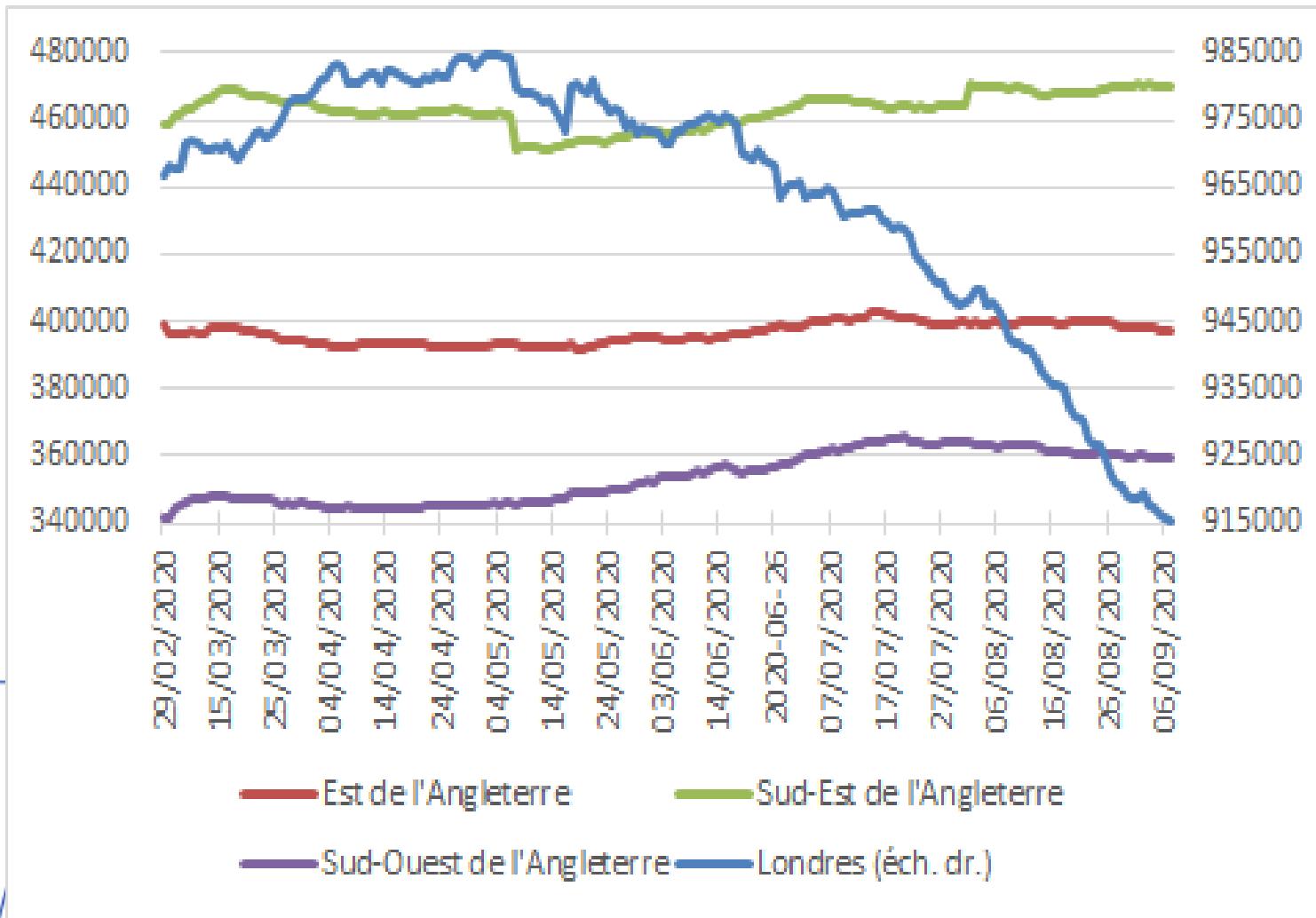
Un marché quasi gelé : offres encore disponibles après un mois sur Internet (à rapprocher d'une chute de 95% des transactions) (en unités)



// Sources : Zoopla, calculs des auteurs.

BANQUE DE FRANCE

Prix de vente moyen affiché par région (en livres sterling par logement)



Sources : Zoopla et calculs des auteurs.



Les statistiques issues du *web-scraping* sont complémentaires des statistiques officielles et permettent un suivi en temps réel et granulaire.

De façon plus structurelle, elles permettent aussi un calcul de prix (ou de loyers) en niveaux.

Les deux sources sont utiles notamment pour les banques centrales : risques de retournement, stabilité financière...



LE WEB-SCRAPING PERMET UN CALCUL DE PRIX EN NIVEAU

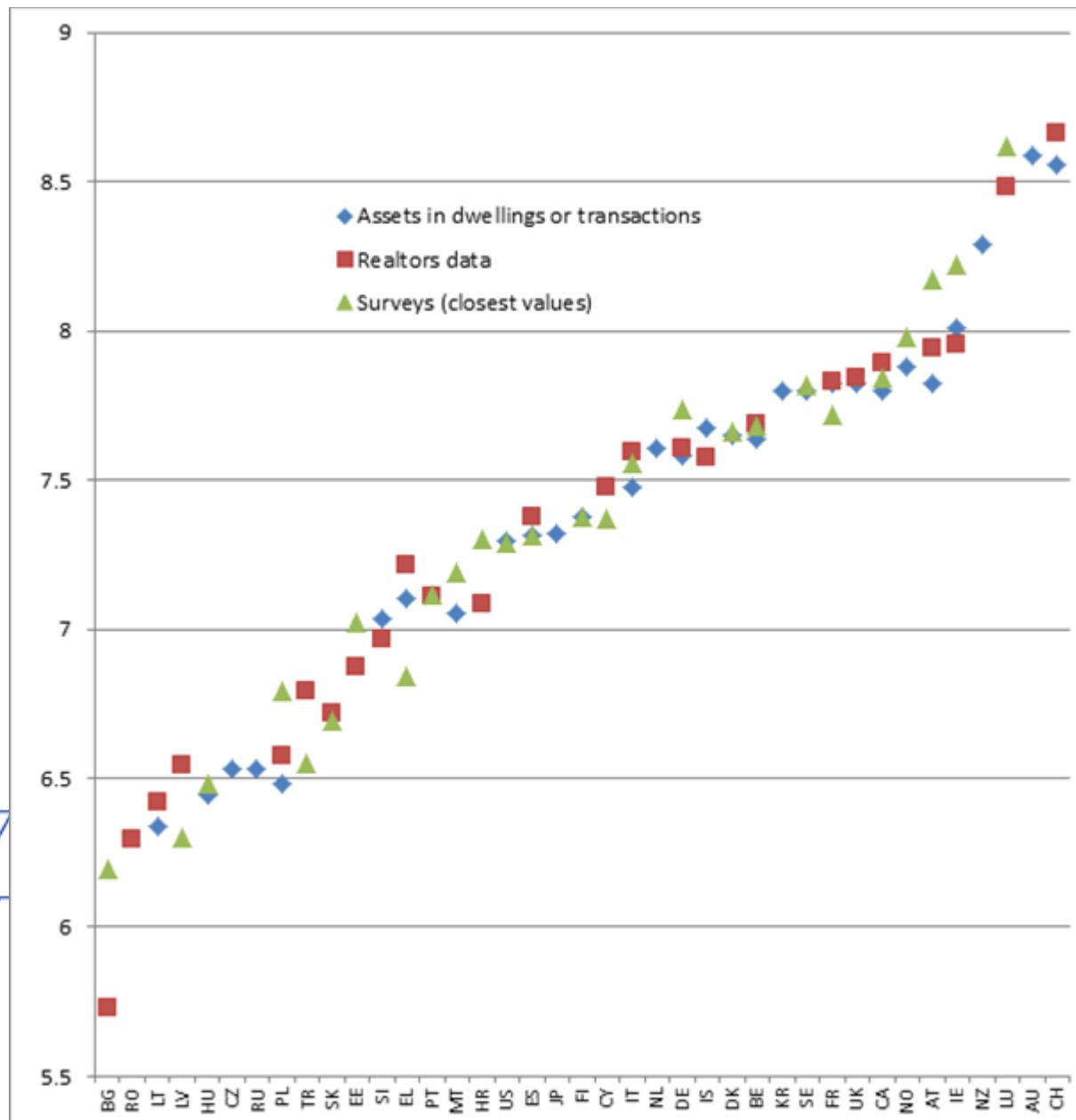
Méthode basée sur données individuelles (*web-scraping*, voire transactions : IE, FR ou MT), au niveau infra-national i (région, municipalité...) avec w_i le stock de m^2 :

$$p = \frac{\sum_{i=1}^n w_i * \text{prix moyen pondéré au } m^2_i}{\sum_{i=1}^n w_i} \quad (2)$$

Avec prix moyen pondéré au m^2_i égal à la somme des prix divisée par la somme des surfaces, équivalent à :

$$\text{prix moyen pondéré au } m^2_i = \frac{\sum_{j=1}^n \text{surface}(j) * \text{prix au } m^2(j)}{\sum_{j=1}^n \text{surface}(j)}$$

WEB-SCRAPING : DES ORDRES DE GRANDEUR COHÉRENTS



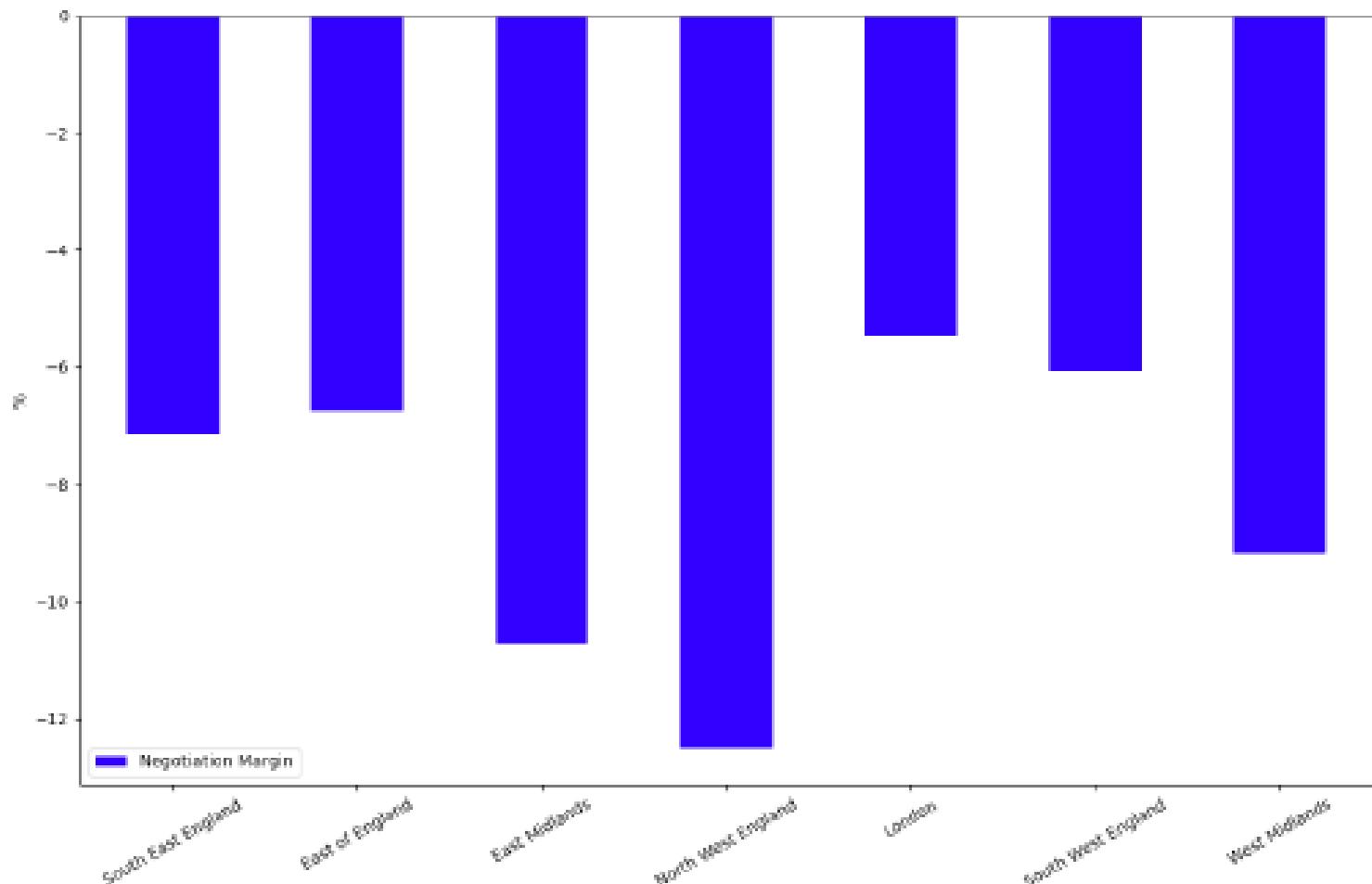
Quelques applications possibles :

- Estimer la richesse immobilière (nationale ou régionale), surtout dans les pays qui n'en publient pas
- Prévoir les retournements (information additionnelle par rapport aux indices) : prix sur revenu, capacité d'achat, marges de négociation si données de transactions sont publiées et les prix postés



Figure 8: Negotiation margin per region

Source: *Real-estate websites, notarial data, and authors' calculation*



Prolongement : le projet Alter Hous en liaison avec le réseau REFINÉ :

- Proposer des statistiques alternatives (non officielles mais présentant des principes communs de fiabilité / représentativité) issues du *web-scraping* ou d'autres sources (satellites...) et des articles méthodologiques sur l'immobilier
- Projet commun avec la Commission européenne, l'OCDE, la Banque mondiale, l'AMSE, le LIEPP, l'ESRI irlandais...





CONCLUSIONS / PERSPECTIVES / ENJEUX

- Approfondir la comparaison avec les statistiques officielles, en gardant en tête les biais possibles (marges de négociation : quels écarts avec les derniers prix postés ?
Travaux en cours)
- Extension aux loyers, au prix des terrains et à l'immobilier commercial (en cours)
- Utilisation de la granularité pour étudier les effets de certaines politiques
- Extension de l'approche à d'autres pays : travaux en cours avec la Banque mondiale
- Questions juridiques liées au *web-craping* et de sécurisation des sources, qui peuvent changer au cours du temps, voire disparaître...



ANNEXES

MOTIVATIONS

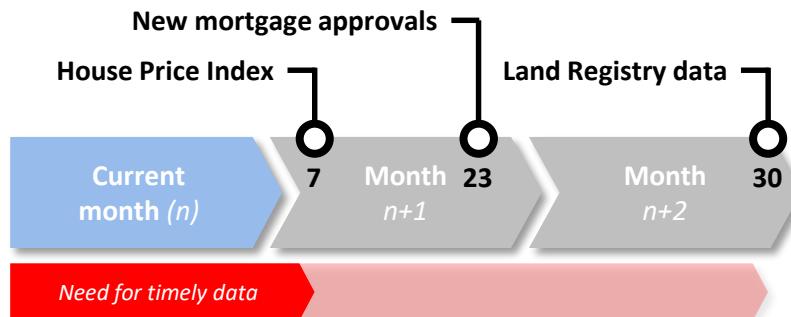
INFORMATION AVAILABLE PUBLICLY AND IN REAL-TIME

Lags in official statistics ...

- Official statistics published only **after month end** and often only aggregated at **national level** or at large region level (Wales, Scotland, Northern Ireland, and England – with potentially a rough decomposition of the latter)

... but information publicly available

- Price information available **publicly** on real-estate website, in **real-time**, and with a high **granularity**
- Evidence that a **large share of offers are posted on the web** (92% of realtors post their ads on those websites)

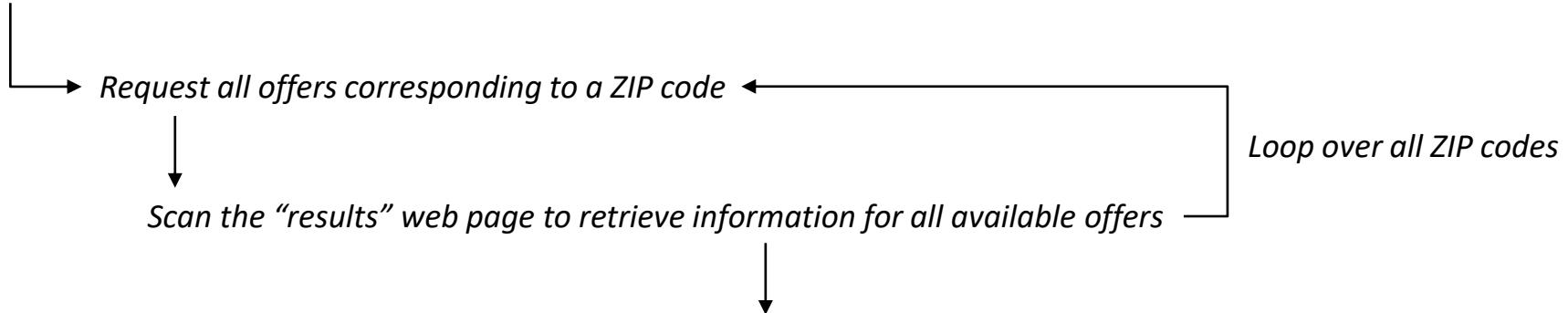


Real-time housing price index
(possibly more granular too)



DATA WEB-SCRAPING

On each of the 5 websites (Zoopla, Rightmove, OnTheMarket, S1homes, and PropertyPal)



ZIP code	City	Address	Surface	# of bedrooms	# of bathrooms	# of living rooms	Transaction type	Building type	Type of good	Price	General description
E1 7A E	London	-	45	2	1	1	Auction	New	Flat	525,000	"Lovely flat with large terrace"
E1 7A E	London	-	130	5	2	1	Sale	Ancient	House	4,500,000	"Garage and safe neighbourhood"

...

E1 7A E	London	-	-	2	2	2	Sale	Ancient	Duplex	850,000	"Balcony and view on Westminster"
---------	--------	---	---	---	---	---	------	---------	--------	---------	-----------------------------------



Around 1.5 million offers in total per day

Figure 13: Median selling prices by region

Source: Zoopla and authors' calculation

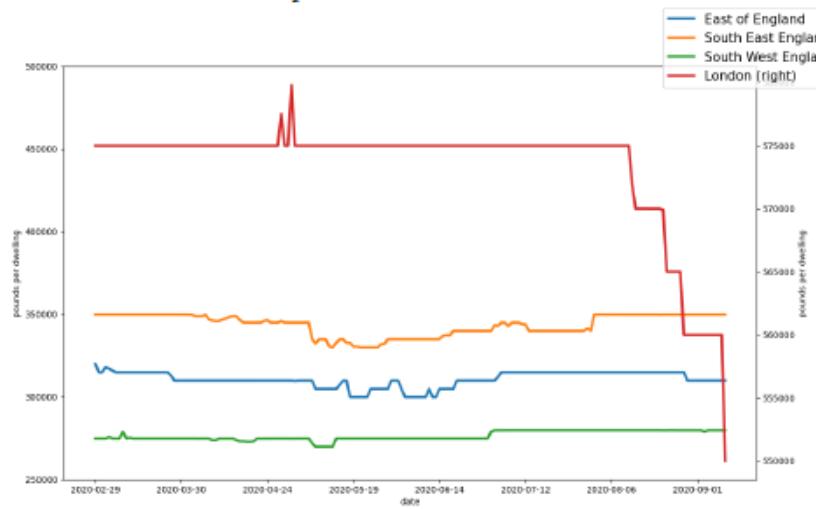
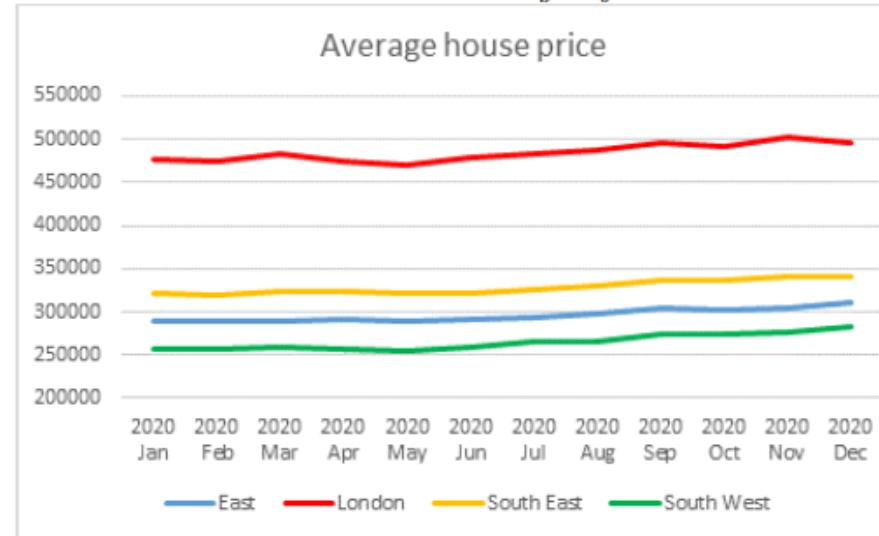


Figure 14: Statistically-adjusted geometrical average selling prices by region

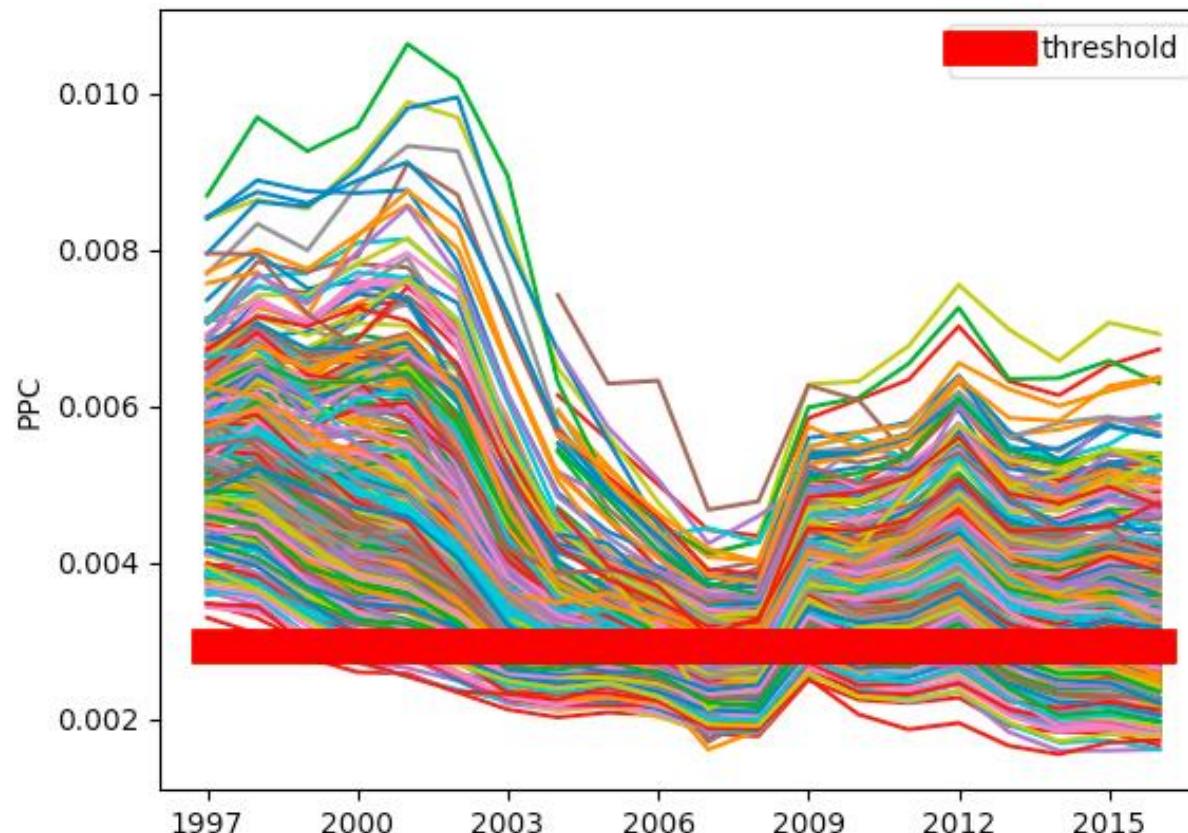
Source: Land Registry



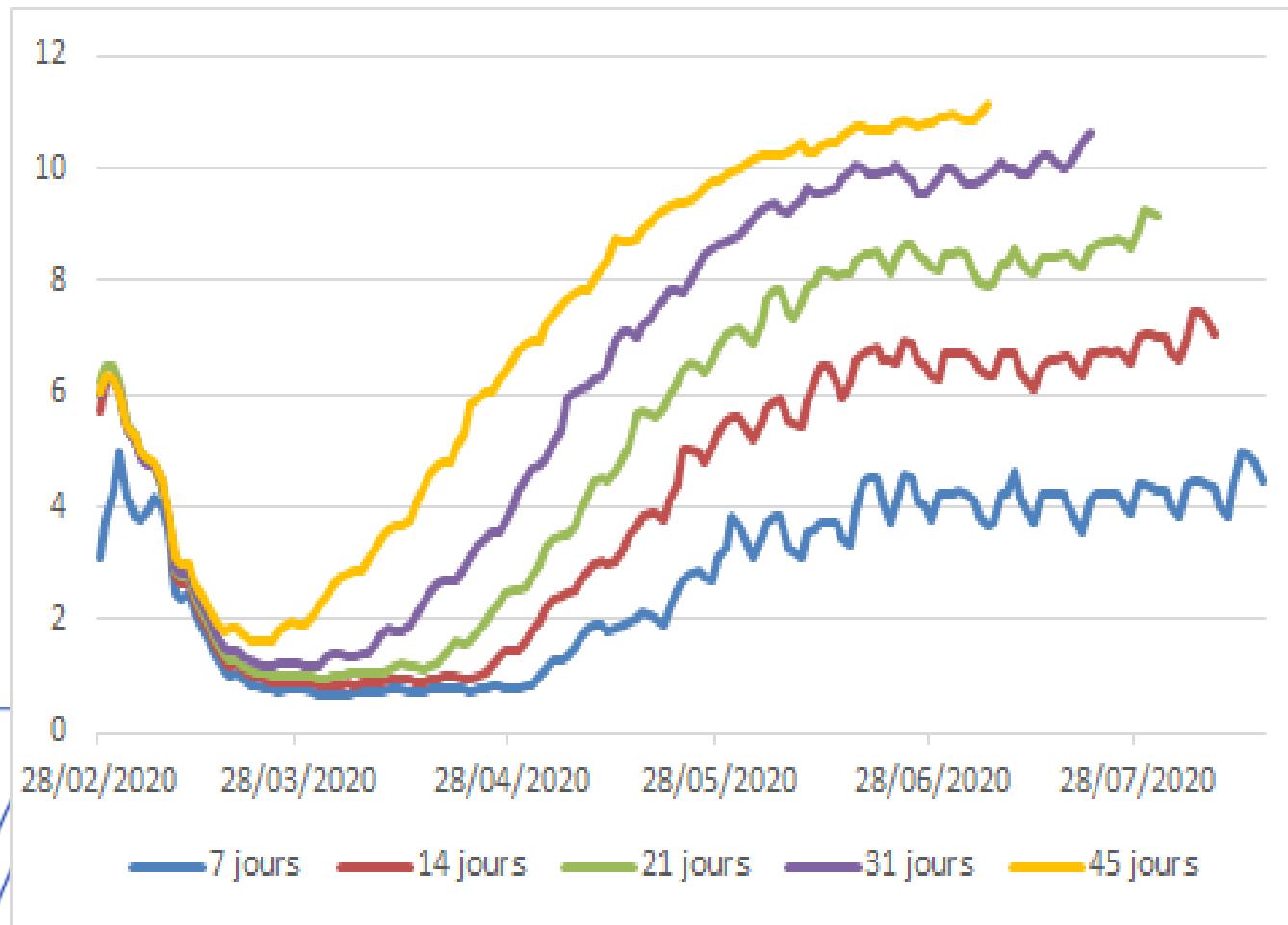
	IS	NO	CH	TR	AU	CA	HK	JP	NZ	RU	KR	US	
1970		129.7						158.2				157.6	
1971		133.3						191.3				162.5	
1972		143.2						234.7				160.6	
1973		165.8						312.0				174.4	
1974		200.2	650.3					340.3				183.5	
1975		226.2	619.4			284.0		323.1				211.4	
1976		272.7	655.8			329.9		354.4				234.0	
1977		277.9	742.0			294.8		417.3				241.5	
1978		266.0	853.9			254.9		486.3				245.2	
1979		269.7	897.1			272.2		410.3				266.5	
1980		311.6	978.3			331.6		602.4				316.1	
1981		368.2	1278.2		747.1	460.4		744.5	307.6			400.5	
1982		371.3	1356.0		754.3	469.4		842.2	402.1			464.1	
1983		424.5	1511.1		854.9	574.3		1045.8	463.7			563.9	
1984		479.7	1566.6		1021.9	628.7		1162.6	443.6			686.5	
1985		553.8	1641.6		742.8	499.7		1190.0	425.6			576.2	
1986		617.7	1844.6		642.1	488.0		1272.2	409.4			510.3	
1987		646.6	2052.8		606.1	493.3		1477.5	492.0			446.1	
1988		709.2	2231.9		977.1	703.4		1690.3	590.9			522.3	
1989		646.2	2404.1		1103.7	801.2		1552.4	580.3			538.7	
1990		589.6	2537.7		960.3	682.9		1639.5	521.4			487.7	
1991		553.7	2390.7		987.9	727.0		1887.2	485.2			502.6	
1992		502.7	2357.1		1005.2	753.3		2012.7	515.1			570.5	
1993		507.2	2384.9		1102.9	797.5		2333.3	631.8			634.2	
1994		578.8	2440.7		1185.5	707.7		2317.5	746.8			594.6	
1995		621.7	2501.9		1077.7	650.8		2063.0	779.9			570.9	
1996		698.8	2119.6		1218.8	678.7		1882.5	977.2			616.4	
1997		777.1	2155.3		1179.6	754.4		1884.2	968.3		933.1	719.5	
1998		789.7	2133.5		1127.4	652.5		2004.1	816.2		1130.2	713.9	
1999		963.9	2134.1		1489.2	828.9	10923.7	2509.6	951.5		1376.4	879.8	
2000		999.7	1096.0	2249.8		1483.2	904.2	10524.0	2320.6	868.2		1354.2	1013.8
2001		912.2	1213.1	2332.1		1601.1	939.4	9760.0	2056.6	880.2		1426.4	1144.5
2002		1036.5	1392.0	2377.0		1770.5	863.4	7290.2	1805.0	1029.9	258.8	1553.7	1029.7
2003		1098.6	1223.7	2247.2		2308.0	953.6	5353.1	1559.7	1278.9	287.4	1398.8	920.6
2004		1302.7	1377.3	2351.7		2359.8	1020.5	6280.4	1416.3	1524.6	359.5	1504.5	934.3
2005		1895.0	1537.9	2419.3		2604.4	1318.1	8574.5	1354.6	1907.6	494.2	1809.9	1191.5
2006		1772.5	1694.8	2486.0		2687.2	1322.4	7717.2	1163.3	1944.2	716.2	1867.1	1130.3
2007		1965.0	1974.8	2555.8		2958.5	1560.2	7688.9	1095.8	2122.5	997.1	1825.4	1011.5
2008		1128.5	1594.6	2955.7		2541.8	1398.6	9528.9	1442.3	1595.6	1081.2	1464.6	983.4
2009		962.9	1909.0	2972.0		3349.5	1527.7	9257.0	1287.0	1917.7	1028.2	1630.2	894.1
2010		1092.0	2198.1	3623.8	748.1	4558.5	1889.2	12389.1	1600.1	2250.8	1177.5	1852.7	935.1
2011		1106.3	2388.2	3964.9	680.0	4605.1	1999.1	15443.3	1736.3	2341.9	911.3	1943.8	925.5
2012		1106.4	2690.9	4186.4	771.9	4594.2	2107.8	17196.1	1517.8	2558.4	1085.6	2097.1	934.4
2013		1252.6	2459.6	4239.8	678.5	4037.5	1935.9	19324.3	1211.1	2670.7	1002.7	2010.5	958.2
2014		1396.4	2336.9	4446.1	796.1	4579.4	2126.3	23256.2	1225.6	3071.0	636.8	2236.4	1144.5
2015		1650.0	2334.6	5062.5	812.7	4970.0	2084.8	29986.0	1390.7	3347.0	578.0	2378.2	1345.6
2016		2148.6	2640.7	5190.0	781.5	5351.5	2442.1	29836.3	1510.0	3977.1	686.7	2438.5	1474.0
2017		2446.5	2559.7	4855.7	704.0	5514.9	2582.9	30381.4	1416.0	613.3	2448.3	1382.3	

Purchasing-power-capacity compared with threshold signalling an housing bubble – by locality

Sources: ONS and authors' calculation



Part des offres dont le prix est révisé, selon le nombre de jours
(en %)



Sources : Zoopla, calculs des auteurs

INDICATORS

RENT-TO-INCOME RATIO

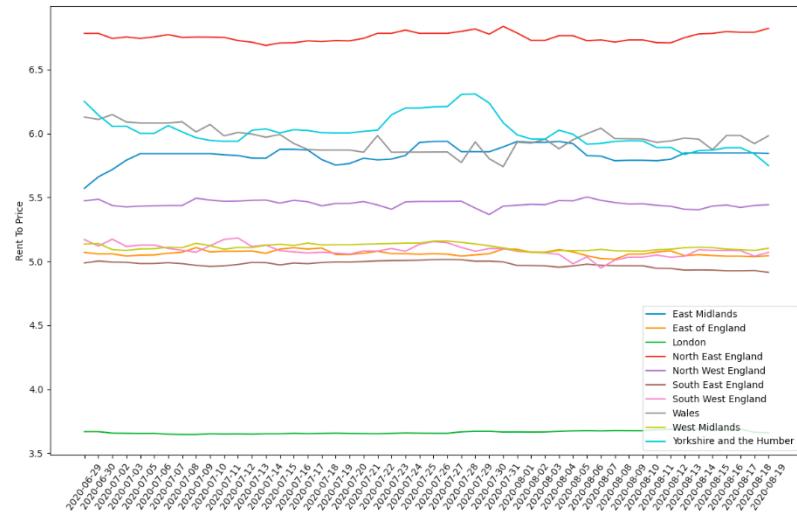
- Compare each **rental offer** with its K closest **sale offers** using a K-nearest neighbor algorithm based on similar characteristics (location, surface, number of bedrooms, number of bathrooms, number of living rooms)

$$RTI = \frac{1}{N} \sum_n \frac{\text{Annual}_\text{Rent}_n}{\frac{1}{K} \sum_k \text{House}_\text{Price}_{n,k}}$$

- We fix K=10 as a result of a trade-off: increasing K results in comparing the rental offer with less and less similar dwellings, but a low K introduces a large variability in the index
- To reflect the unequal coverage of our data, the ratio is estimated at local level then aggregated at national level **weighing each regional index by the stock of dwellings in that region**
(source: ONS)

Rent-to-price ratios – by region

Sources: Zoopla and authors' calculation



- Little variation during the crisis but **large discrepancies across regions**. In particular, very low ratio in London – might signal overvaluation to some extent (in line with Marsden, 2015 or Petris et al., *in press*)
- At national level, index around 5.1%

PRICE-TO-INCOME RATIO & PURCHASING-POWER- CAPACITY

- Combine our dataset with **ONS' data on disposable income** – also available to some extent at local level – to derive a price-to-income ratio = house price divided by disposable income
- However, as defined in the literature, price-to-income ratio does not consider the **evolution of interest rates**. Propose an alternative index (**the purchasing-power-capacity**) computed as:

$$PPC = \frac{\omega}{r_t} \cdot PTI_t \cdot \left(1 - \frac{1}{(1 + r_t)^d} \right)$$

ω = maximum effort rate (33%)

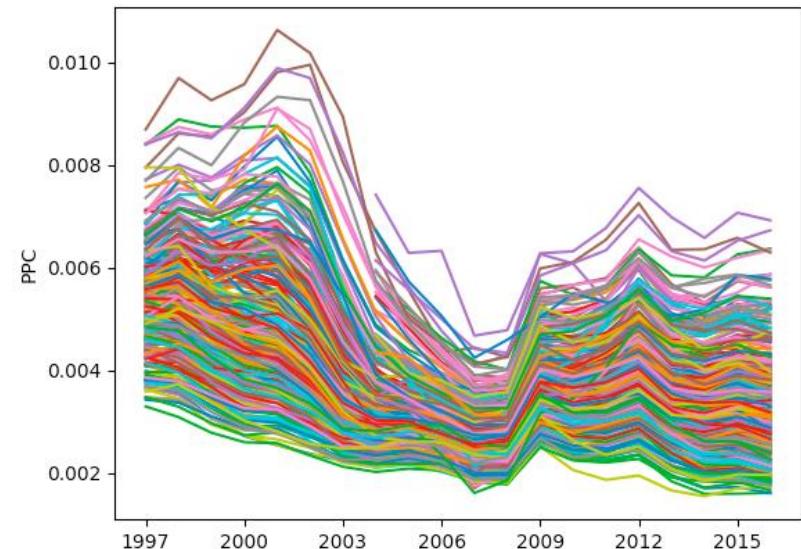
r_t = mortgage rate (source: BoE)

PTI_t = price-to-income ratio

d = median duration (15 years)

Purchasing-power-capacity – by locality

Sources: ONS and authors' calculation



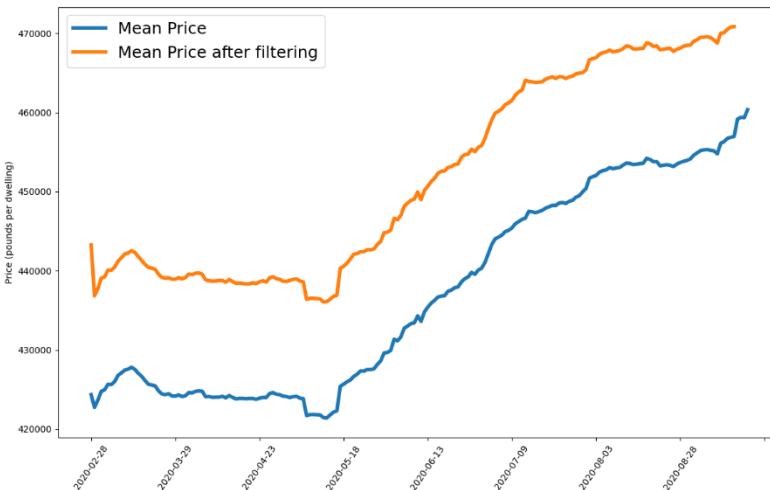
- Large discrepancies but historical evidence for a **reversal in trend after the GFC**
- Empirically, find a **higher signalling power** (share of true detection of a forthcoming housing crisis) for the PPC than for the price-to-income ratio generally used in the literature

- ① Ensure **consistency** across units (price in GBP, area in squared meters, rent as monthly amount), numeric (surface/price can given be as a range), and text data (lowering characters, lemmatisation) in order to enhance comparability across offers
- ② Remove **duplicates** – which can arise even on the same website
- ③ Exclude **commercial real-estate, miscellaneous offers** (garage, land, mobile home, bungalow, etc.) and **auctions** (since the price displayed would be the reserve price)
- ④ Remove **outliers** by winsorising at 1%
- ⑤ Using natural language processing on the general description of the offer, add **dummies to account for the presence of additional facilities** (e.g. presence of a garage, a garden, or a terrace) that might a premium to the price

REAL-TIME MONITORING UK DURING COVID-19 (3/4)

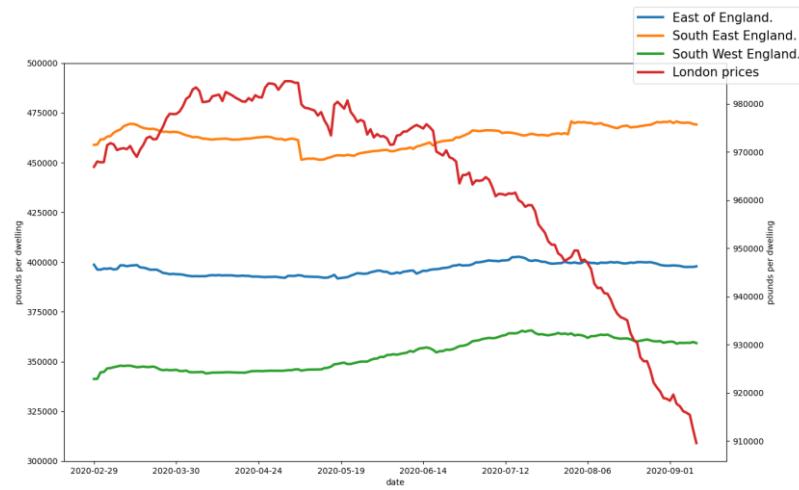
Selling price evolution – all UK

Sources: Zoopla and authors' calculation



Selling price evolution – by region

Sources: Zoopla and authors' calculation



- At national level: **slight decrease** from mid-February to the beginning of May (first lockdown), followed by a steady increase afterwards – with no apparent effect of the second lockdown
- At local level: **steady decline in London since the end of the lockdown** – in contrast with an increase in other regions (East England, South- East England, and South-West England).
- Looking at evolutions by quantiles shows that these **patterns are shared across all categories of dwellings**

Table 5. Signalling: possible cases

Crisis signal No crisis signal	No crisis episodes (NCE)				Crisis episodes (CE)			
	False alert (FA)		True negative signal		True positive signal		Missed crisis (MC)	
Threshold	Signal power	% False alerts (type 1 error)	% Missed crises (type 2 error)	# Crisis years	Size of the sample	2SD lowerbound	2SD upperbound	
(1)	(3)	(4)	(5)	(6)	(7)	(8)	(9)	
Full sample								
PTI in level	10.21	0.32	0.28	0.4	40	722	9.93	10.48
Max PTI over 3 previous years	10.95	0.44	0.21	0.35	34	680	9.84	12.06
PTI/whole period mean PTI	1.03	0.26	0.39	0.35	40	722	0.93	1.13
PTI/real time mean PTI	1.20	0.06	0.19	0.75	40	722	0.87	1.53
Sample where the PTI of the first five years of each country are deleted								
PTI in level	10.20	0.33	0.27	0.4	30	593	9.94	10.47
Max PTI over 3 previous years	10.95	0.49	0.18	0.33	27	538	10.44	11.46
PTI/whole period mean PTI	0.96	0.22	0.55	0.23	30	593	0.81	1.11
PTI /real time mean PTI	1.33	0.07	0.10	0.83	30	593	0.96	1.70
Sample where the PTI of the first ten years of each country are deleted								
PTI in level	10.19	0.39	0.24	0.37	19	452	9.93	10.44
Max PTI over 3 previous years	10.02	0.49	0.26	0.25	16	408	9.06	10.98
PTI/whole period mean PTI	1.14	0.36	0.22	0.42	19	452	1.03	1.25
PTI /real time mean PTI	1.25	0.08	0.18	0.74	19	452	0.94	1.55