

Universität Leipzig
Institut für Informatik



Kategorisierung von Produktangeboten in großen Katalogen

Hanna Köpcke
WDI-Lab

28.10.2010

Produktkatalog

- Systematisch geordnete Sammlung von Informationen zu Produkten oder Dienstleistungen
- Klassifikation anhand eines gemeinsamen Merkmales, z.B. Hersteller, Nutzung, Preis, Farbe, Feature
- Beispiele:



Kategorien		
Antiquitäten & Kunst	Filme & DVDs	Musikinstrumente
Audio & Hi-Fi	Foto & Camcorder	PC- & Videospiele
Auto & Motorrad: Fahrzeuge	Garten	Reise
Auto & Motorrad: Teile	Handy & Organizer	Sammeln & Seltenes
Baby	Haushaltsgeräte	Software
Beauty & Gesundheit	Heimwerker	Spielzeug
Briefmarken	Immobilien	Sport
Bücher	Kleidung & Accessoires	Tickets
Büro & Schreibwaren	Modellbau	Tierwelt
Business & Industrie	Möbel & Wohnen	TV, Video & Elektronik
Computer	Münzen	Uhren & Schmuck
Feinschmecker	Musik	



Alle Kategorien ansehen	
Bücher	>
Musik, DVD & Games	>
Computer & Büro	> Notebooks & PCs
Elektronik & Foto	> PC-Zubehör & Monitore
Küche, Haus & Garten	> PC-Komponenten
Baumarkt & Auto	> Software
Lebensmittel & Drogerie	> PC- & Video-Games
Spielzeug & Baby	> Drucker & Tintenpatronen
Kleidung, Schuhe & Uhren	> Bürobedarf
Sport & Freizeit	>



Computer • Notebooks/Laptops • PCs • Server • Workstations • Netbooks	Netzwerke • Netzwerk-Switches • Netzwerk/Internet-adapter • Router • Telefone	Audio u. Video • MP3-Player u. -Recorder • DVD-Player/-Recorder • Lautsprecher • Computerlautsprecher • Medienplayer/-recorder • TV-Tuner-Karten • Home-Kino Systeme	Küche u. Haushalt • Küchenherde & Kocher • Küchschneide • Waschmaschinen • Staubsauger
Komponenten • PC-Speicher/RAM • Prozessoren • Grafikkarten • Motherboards • Audokarten	Datenspeicher • Festplatten / HDD • Flash-Speicher • CD/DVD-Brenner • Kartenleser • CD/DVD-R/DW Laufwerke	PDA, GPS u. Handy • PDAs • Handys • GPS-Navigationsysteme	Verbrauchsmaterialien u. Zubehör für Büromaschinen • Leimröhren • Akkumulatoren • Papier- und Tintenpatronen • Papier-Perforatoren u. Zubehör • Bindesysteme
Drucken und Scannen • Multifunktionsgeräte • Laserdrucker • CD-Label-Drucker • Tintenstrahl-Drucker • Fotodrucker	Monitore, TV u. Displaygeräte • Flachbildschirme / TFTs • LCD TVs • Plasma-Bildschirme	Kameras • Digitale Kameras • Digitale Videokameras • Webcams	Körperpflege • Rasiergeräte der Männer • Haartrockner/Föhne • elektrische Zahnbürsten • Solar- / Sonnenstudio
Software • Archivier- & Sicherheits-Software • Desktop-Publishing-Software • Betriebssysteme • Software-Güter • Navigations-Software	Taschen/Etuis • Notebooktaschen • Kameraschalen und Rückkäcke • Kuckucke	Einbaugeräte • Tastaturen • Mäuse • Speicherkontroll	

Kategorisierung

- Einordnung von Produktangeboten in eine (oder mehrere) Kategorien des Katalogs
- Schnelleres Auffinden der gesuchten Information
- Herausforderungen:
 - Große Anzahl an Kategorien und Angeboten
 - Große Unterschiede zwischen Katalogen
 - Struktur, Bezeichnungen
 - Falsche Zuordnungen
 - Mehrere mögliche Kategorisierungen

(Semi-)automatische Kategorisierung

- 2-stufiger Prozess:
 1. Training
 - Voraussetzung: Angebote für die die Kategoriezuordnung bekannt ist
 2. Kategorisierung
 - Anwendung auf:
 - Kompletten Bestand eines Shops
 - Neue Angebote
 - Aktualisierte Angebote

Trainingsphase

- Tokenisierung
- Entfernung von Stoppwörtern
- Modifizierter Naive Bayes Ansatz
- Verbesserung der Trefferquote durch inkrementelles Lernen

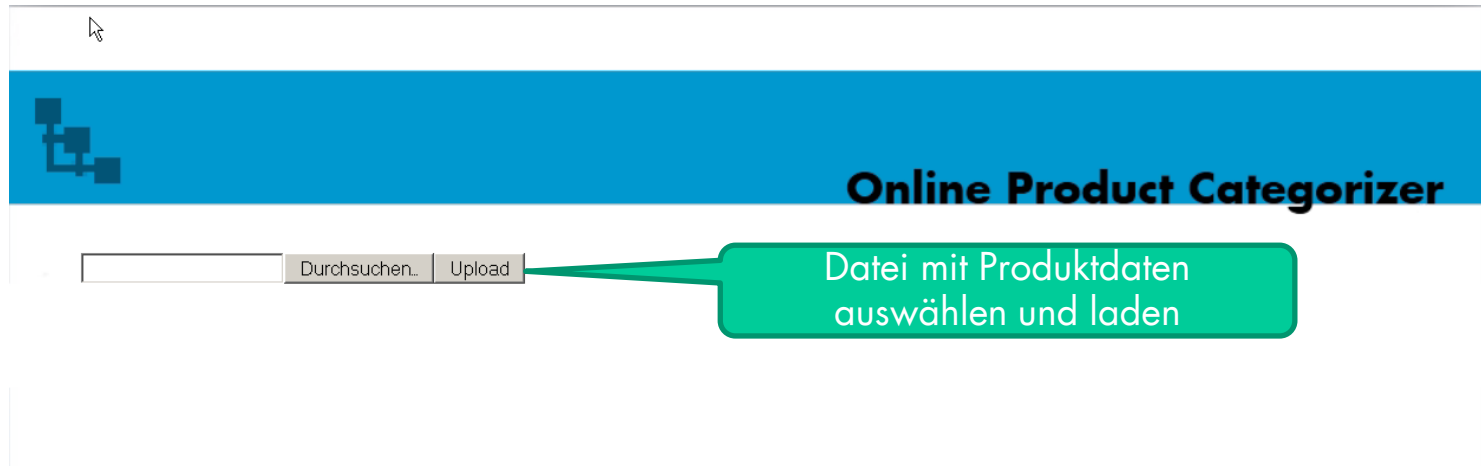
a4 amp bis blatt brother color dpi druck- druckauflosung
drucken drucker druckgeschwindigkeit duplexdruck farbe faxen
funktionen fur gewicht gt hp jahre kopier- kopieren kopierformat
kopiergeschwindigkeit laserdrucker laserjet max mb min mit
multifunktionsdrucker multifunktionsgerate multifunktionskopierer
netzwerkfähig oki papiervorrat samsung scan-auflosung scanauflosung scannen
schnittstelle seiten sharp speicher SW usb von vorlageneinzug zu

angabe anmerkungen anschlusswerte bitte bleiben braucht ca cm elektrogerate energie
energieeffizienzklasse energieverbrauch fach-beratung farbe
gefrierkapazität gefrierteil gerate geratema gerne gro
helfen hierfür jahresenergieverbrauch kg klimaklasse kuhlgerate kuhlteil
kwh lagerzeit liter lohnt maximale modelle nutzinhalt persönlich
regel richtigen schnell seitenteile sparsame std störung technik-experten top-feature
turanschlag vorteile wahl wechselbar wei wenden

Vorteile

- Reduzierung des manuellen Aufwandes
- Überprüfung bestehender Kategorisierungen
- Unabhängig von Änderungen der Quellkataloge
- Anpassbar an eigenen Katalog
- Skalierbar für große Anzahl von Kategorien und Angeboten

Online Product Categorizer



Starten der Kategorisierung

E:\workspace_hk\Product Durchsuchen... Upload **Categorise**

Anzeige der Produkte

Anzeige des Katalogs

Category tree		description	price
[-] Foto-, Videokamera & Co.	1	Kino und Hermes reisen durch verschiedene Länder, in denen sie ihnen...	7.5
[-] Computer & Software	2	Grundkurs Marmorieren Anne Pieper;	8.95
[-] Schnäppchen & Gebrauchtware	3	"Schmidt - Kniffel ""Kniffel-Block"" , 3er-Set, für 1440 Spiele "	6.49
[-] Audio, Video & TV	4	Hört mal, was da klinglt!, m. Audio-CD Conny Frühauf;Christine Werner;	28.9
[-] Buch, Hörbuch & Kalender	5	Aktmalerei, m. DVD Martin Thomas;	19.9
[-] Beauty, Wellness & Drogerieartikel	6	Das große Buch der Gouachemalerei Ute Schmidt;	16.8
[-] Downloads zum Verkauf & Verleih	7	Mein erstes Bildlexikon Haustiere	9.95
[-] Telekommunikation	8	Fenstersticker aus Windowcolor Stefanie Thomas;Heike Roland;	7.5
[-] Game & Konsole	9	Kinos Reise Keichi Sigsawa;	7.5
[-] Hobby & Spiel	10	Chibi-Manga zeichnen und malen Christopher Hart;	14.95
[-] Haushalt & Wohnen	11	Mein erstes Bildlexikon, Tiere im Wald Antje Kleinlüdern-Depping;	9.95
[-] Sport & Freizeit	12	Die Klarinette, m. Audio-CD Rudolf Mauz;	19.95
[-] Mode & Accessoires	13	Das Grüffelo-Puzzle-Buch Julia Donaldson;Axel Scheffler;	15.9
[-] weitere Produkte	14	Eine kulinarische Entdeckungsreise durch den Chiemgau Petra Wagner;	29.9
	15	Eine kulinarische Entdeckungsreise durch Münsterland und Osnabrück...	29.9
	16	Die Einkreisung Caleb Carr;	9.95
	17	Volevo i pantaloni Lara Cardella;	8.95
	18	Die Unruhezone Jonathan Franzen;	19.9
	19	Asterix, English editionPt.31 Asterix and the Actress; Asterix und Latra...	11.2
	20	Arcodamore Andrea De Carlo;	9.9
	21	Xenophobe's Guide to the Americans Stephanie Faul;	6.5
	22	Winnetou, Engl. ed. Karl May;	34.9
	23	The Study of Language George Yule;	21.3
	24	Der Junge, der Anne Frank liebte Ellen Feldman;	9.95
	25	Die Wirklichkeit der tropischen Mythen Fritz J. Raddatz;	12.5

Kategorisierung

id	title	description	price	category ^
1451	Tiger & Dragon Wang Du Lu;	400 Jahre alt ist das wertvolle Schwert, das den beiden...		5 Action-DVD
1455	Godzilla Terry Rossio;Dean Devlin;Roland Emmerich;Ted E...	In Manhattan herrscht Krieg. Gegen den furchterregendst...		99 Action-DVD
1642	Telekommunikation - Schnurloses Telefon »Gigaset A580...	Schnurloses Telefon »Gigaset A580«, Telefonbuch für bi...		Analog-Telefon
1643	Telekommunikation - Schnurgebundenes Telefon »Eurose...	Schnurgebundenes Telefon »Euroset 5035«, mit Anrufbe...	49.9	Analog-Telefon
1644	Telekommunikation - Schnurloses Telefon »Gigaset C385...	Schnurloses Telefon »Gigaset C385 Duo«, mit digitalem A...	98.76	Analog-Telefon
1645	Telekommunikation - Schnurloses Telefon »Gigaset C380...	Schnurloses Telefon »Gigaset C380 Duo«, mit zusätzlich...	83.29	Analog-Telefon
1646	Telekommunikation - Schnurloses Telefon »Gigaset A585...	Schnurloses Telefon »Gigaset A585«, mit digitalem Anruf...	49.97	Analog-Telefon
1647	Telekommunikation - Zusatz-Mobilteil von Philips,	Zusatz-Mobilteil, für Magic5-Eco-Voice, mit hintergrundbel...	47.59	Analog-Telefon
1648	Telekommunikation - Mobilteil »Gigaset C47H« von Siemens,	Mobilteil »Gigaset C47H« (Mobilteil »Gigaset C47H«)	67.82	Analog-Telefon
1649	Telekommunikation - Schnurloses Telefon »Gigaset E360«...	Schnurloses Telefon »Gigaset E360«, mit großen Tasten,...	77.34	Analog-Telefon
1650	Telekommunikation - Mobilteil »Gigaset S68H« von Siemens,	Mobilteil »Gigaset S68H«, mit Anrufanzeige mit Bild, Farbdi...	71.39	Analog-Telefon
1651	Telekommunikation - Telefon »Euroset 5020« von Siemens,	Schnurgebundenes Telefon »Euroset 5020«, Display: 1-z...	40.45	Analog-Telefon
1652	Easy CA 22 T Home		43.99	Analog-Telefon
1890	DVB-T Antenne, aktiv	DVB-T Antenne, aktiv- aktive Stabantenne für digitalen T...	14.79	Antenne
900	ENDPOINT PROTECTION 11.0 GR	5USER BNDL BUSINESS PACK GR	250	Antiviren-
901	Norton Smartphone Security	5.0/ deutsch/ Vollversion/ DVD Case/ Datenträger: CD	22.4	Antiviren-
902	ANTIVIR PERSONAL EDITION D	PREMIUM 1 YEAR	21.8	Antiviren-
903	Total Protection for SMB Int	25 USER inkl. 1yr Gold Support	647.87	Antiviren-
904	Total Protection for SMB Int	5 USER inkl. 1yr Gold Support	158.06	Antiviren-
905	Total Protection for SMB Int	10 USER inkl. 1yr Gold Support	285.51	Antiviren-
906	Update CA Threat Manager r8.1	25 User - Upgrade from eTrust Antivirus, eTrust PestPatr...	677.92	Antiviren-
907	Total Protection SMB Adv 10U	BOX - 10 USER inkl. 1Jahr Gold Support/ Schutz für: Desk...	534.16	Antiviren-
908	Update CA Threat Manager r8.1	5 User - Upgrade from eTrust Antivirus, eTrust PestPatrol...	135.6	Antiviren-
909	Total Protection SMB Adv 25U	BOX - 25 USER inkl. 1Jahr Gold Support/ Schutz für: Desk...	1218.89	Antiviren-
1990	RAM DDR2-667 2GB/CL5/KIT-2x1024MB	DDR2 2GB / non-ECC / CL5 / 667MHz / ValueRAM / KIT-2x...	60.98	Arbeitsspeicher

Einschränken der Anzeige

Category tree		description	price	category ^
+	Foto-, Videokamera & Co.	...eine Tierpension	" Auf dem PC bereits höchst erfolgreich, halten die possie...	19.99 Nintendo DS-Game
+	Computer & Software	Naruto: Ninja Destiny	Überraschend, blitzschnell und leise schlagen sie zu. Die...	35.99 Nintendo DS-Game
	Schnäppchen & Gebrauchtware	Mario Party DS	4 Spieler können sich gegenseitig beweisen, wer der Be...	39.99 Nintendo DS-Game
+	Audio, Video & TV	- Englisch Buddy	" Buddy macht fit in Englisch. Mit über 5.000 Wörtern, übe...	29.99 Nintendo DS-Game
+	Buch, Hörbuch & Kalender	Aquarium by DS Jowood;		24.99 Nintendo DS-Game
+	Beauty, Wellness & Drogerieartikel	47 Matchstick Puzzle by DS Jowood;		24.99 Nintendo DS-Game
+	Downloads zum Verkauf & Verleih	748 Think Logik Trainer	"Aufgaben in den Bereichen Sprachverständnis, Gedächt...	29.99 Nintendo DS-Game
+	Telekommunikation	749 Sarah, die Hüterin des Einhorns, Nintendo DS-Spiel	Sarah lebt mit ihrem Onkel auf einer Farm am Waldrand. A...	39.99 Nintendo DS-Game
-	Game & Konsole	870 Titanic Mystery	Spannende Aufgaben und knifflige Rätsel Die Titanic stich...	9.9 Nintendo DS-Game
	Konsolen-Game	724 Guitar Hero Greatest Hits	"KOMPLETTE TRACKLISTE VON GUITAR HERO GREATES...	64.99 Nintendo Wii-Game
	Gameboy-Game	727 Super Mario Galaxy	Prinzessin Peach wurde ins Weltall entführt - doch Mario...	49.99 Nintendo Wii-Game
	Gamecube-Game	728 Harvest Moon Magical Melody	Harvest Moon Magical Melody: Ab auf die Farm!	39.99 Nintendo Wii-Game
	Nintendo 64-Game	729 "Wii-Spiel Ban Dai ""Wii Family Trainer: Double Challenge""	Plattform: Nintendo Wii. Genre: Family Entertainment. USK: ...	59.99 Nintendo Wii-Game
	Nintendo DS-Game	732 Mein Fitness Coach - Cardio Workout (Nintendo Wii)	Trainieren Sie mit Ihrem eigenen Coach Fit werden mit der...	19.9 Nintendo Wii-Game
	Nintendo Wii-Game	733 Guinness World Records - The Videogame	Guinness World Records: Das Videospiel steckt voller un...	9.99 Nintendo Wii-Game
	PS2-Game	734 Celebrity Sports Showdown	"Ohne dieses Spiel für die Wii kommt ab sofort keine Party...	39.99 Nintendo Wii-Game
	PS3-Game	735 "Wii Spiel SEGA ""Mario & Sonic bei den Olympischen Win...	Plattform: Nintendo Wii. Genre: Sport/Party. Spieler: 1-4. O...	59.99 Nintendo Wii-Game
	PSP-Game	743 Spongebobs Atlantisches Abenteuer	"Begleiten Sie SpongeBob und den Rest seiner Kumpels a...	19.99 Nintendo Wii-Game
	Sega-Game	751 Speed Racer: Das Videogame		19.99 Nintendo Wii-Game
	Xbox 360-Game	787 Sing it: High School Musical	"Disney Sing It ist ein neues videobasiertes Karaoke-Spiel...	19.99 Nintendo Wii-Game
	Xbox-Game	788 Spongebobs Atlantisches Abenteuer	Ein antikes Medaillon bringt SpongeBob und seine Freund...	9.99 Nintendo Wii-Game
	sonstige Games	885 Guitar Hero Greatest Hits	"KOMPLETTE TRACKLISTE VON GUITAR HERO GREATES...	70.99 Nintendo Wii-Game