

PETER GLENN . COM

a scraping project
by: Ido Farhi





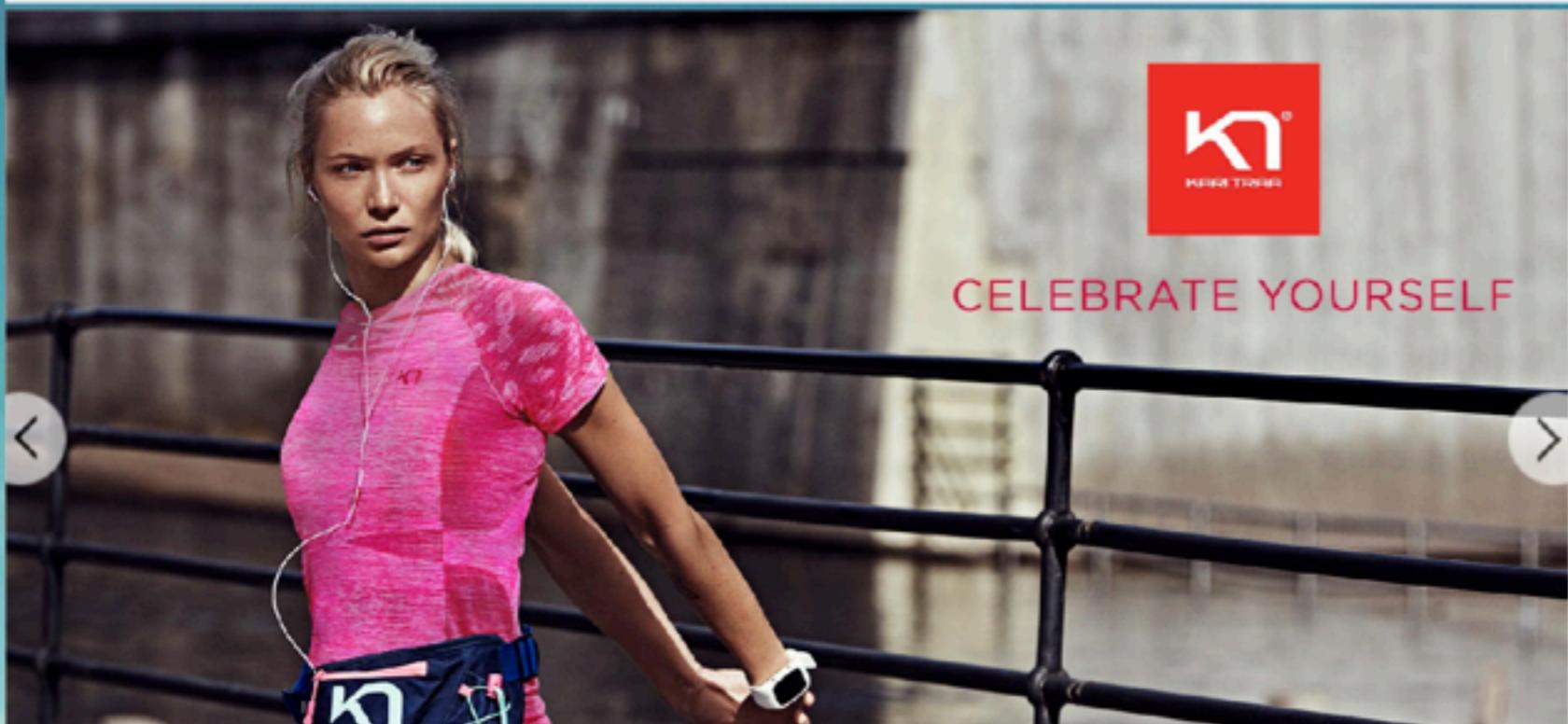
SEARCH by keyword, product name or item #



ORDERS OVER \$49 SHIP FREE

1-800-818-0946

MEN WOMEN KIDS SKI SNOWBOARD RUN TRAVEL WATER SKATES BRANDS SALE



CHECK OUT THE NEW SPRING ARRIVALS ▶





SEARCH by keyword, product name or item #



1-800-818-0946



ORDERS OVER \$49 SHIP FREE

MEN WOMEN KIDS SKI SNOWBOARD RUN TRAVEL WATER SKATES BRANDS SALE

SKATE
to SKI

SHOP THE SKATE-TO-SKI CROSS TRAINING SERIES ▾

BRANDS

A

Abrams Books
Adidas
AFRC
Ahnu
Airblaster
Airhole
Allsport
Alp-n-Rock Clothing
Alpinestars
Altra
Altro Optics
Amphipod
Anita
Anon
Apex Ski Boots
ARBOR
Arc'teryx
Arcade
Armada
Astis
Astral
Atomic

B

Balega
Bearpaw
Bent Metal
Bergans of Norway
Black Diamond Equipment
Blizzard
Blundstone
Bogner
Bomber Ski
BOS. & CO.
Boulder Gear
Brooks
Brunton
Bugatchi
Bula
Burton

C

Camaro
Camas Designs
Camelbak
Capita
CEP
Coal
Colmar
Columbia
CP
Crab Grab
crazeHeads
CWB
CWX

D

DAKINE
Dalbello
DC Shoes
Descente
DOUBLE DIAMOND
Dragon Goggles
Dynastar

E

Eisbar
Elan
Elan Blanc
Electric
Erin Snow
Euro Socks

F

Feeutes!
FERA
Fire + Ice
FISCHER
Fitletic
Fjallraven

G

GIRO
Globe
GNU Snowboards
GOLDBERGH
Goldwin
GoPole

H

Harmony
Head
Headsweats
Helly Hansen
Hi-Tec
High Range Gear



Information scraped from each page

The North Face Condor Triclimate Ski Jacket (Men's)

Was \$290.00
SALE: \$179.99

FREE SHIPPING

4.3 (94 reviews) [Read 94 Reviews](#) [Write a Review](#)

Quick Specs

Fit: Relaxed [?](#)
Waterproofing: 8/10 [?](#)

Choose Size:

S M L XL XXL

Choose Color:

ADD TO CART

ADD TO LIST

Starting at \$16 a month with [Affirm](#) [Learn More »](#)

360° VIDEO

FIND in a store

Outfit Builder

ADD to builder

PRODUCT DETAILS



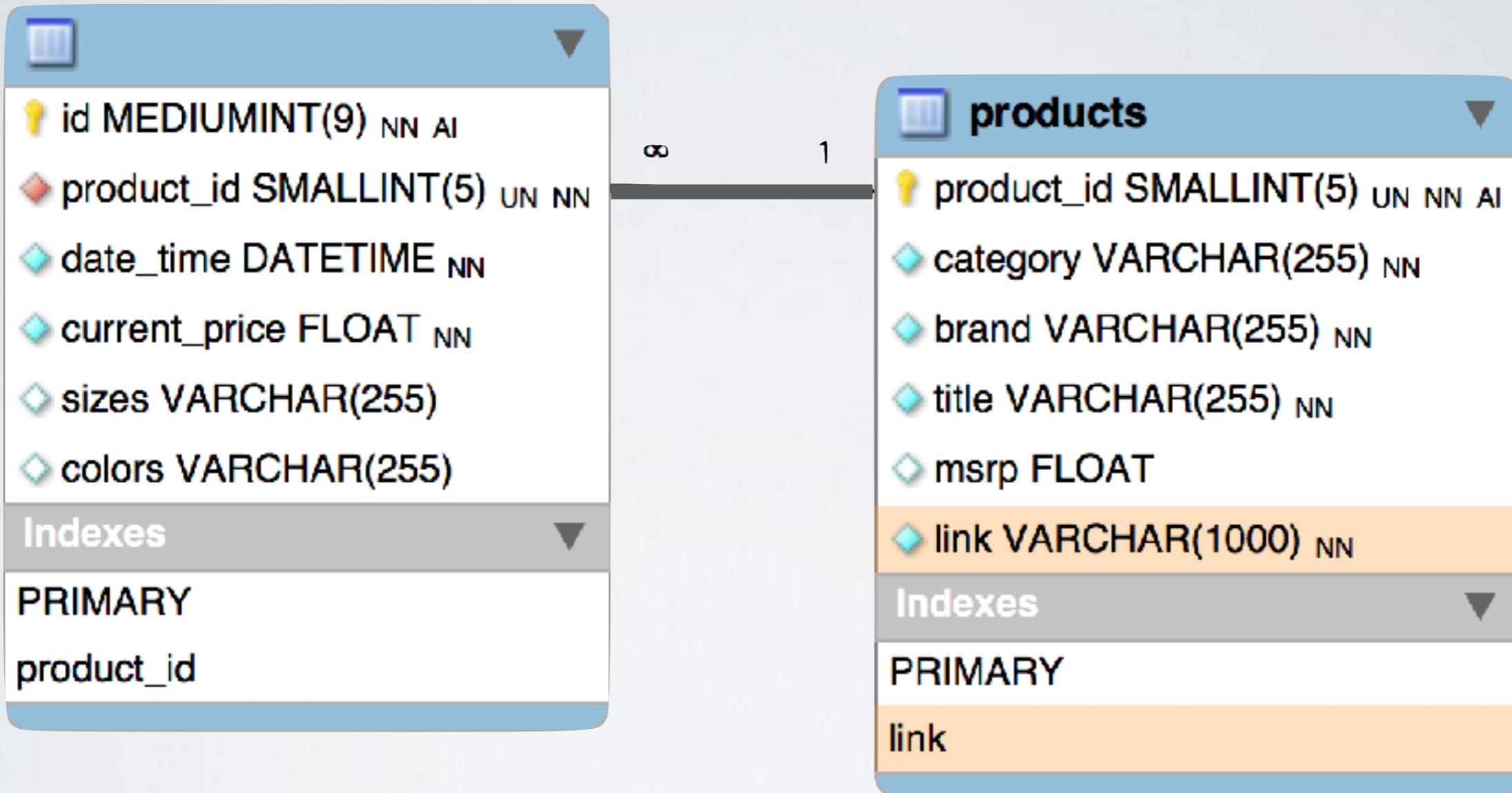
Like 2

FIT

WATERPROOFING

MATCHING ITEMS

DATABASE (mySQL)



- Create a category table
- Create a brand table

SOME NUMBERS

- Distinct items scraped: ~5,500
- Rows in database: ~100,000
- Time to scrape each item: 2 seconds
- Total run time: ~3 hours ... not reasonable

MULTIPROCESSING

VS.

MULTITHREADING

- **Process:** running code.
- A process can have multiple threads.
They run the same code belonging to the parent process in parallel.
- Threads require fewer resources and are “cheaper” to create and run

MULTITHREADING

- Based on `concurrent.futures`
- Part of the standard library that provides a higher-level control over threads.
- The threads are modelled as asynchronous tasks.
- `concurrent.futures` => `threading` => `_thread`
- Each thread collects one brand at a time
- Total run time: ~20 minutes (~10 on AWS)

THE “\$ > 1,000” BUG

Home > Ski > Skis > Blizzard Quattro 8.4 TI Ski System with Bindings (Men's)

Blizzard Quattro 8.4 TI Ski System with Bindings (Men's)



Was \$1,080.00
SALE: \$719.99

FREE SHIPPING

0.0 (No reviews) Q&A Ask a Question

Be the first to Write a Review

Quick Specs

Profile: Rocker ?

Waist Width: 84 ?

More Info

Choose Size:

167 174 181

ADD TO CART

ADD TO LIST

Starting at \$64 a month with **affirm** Learn More »

 FIND in a store

 Outfit Builder

ADD to builder

A red arrow points from the word "SALE" in the price text to the "33% Off!" badge.

```
295     try:  
296         # Catch exception in case no msrp found  
297         msrp = float(soup.find('span', itemprop='price_msrp').text[1:])  
298     except AttributeError:  
299         # no different price & msrp. They are equal.  
300         msrp = current_price
```

- Search for item msrp, if not found, set it equal to the current price.
- Can you guess what happens when msrp is \$1,000 or more ?
- How is multithreading related ?
- Then how was this bug caught ???
- Moral of the story...

CHALLENGES

- How to scrape the whole site ?
- Long I/O times
- Server uptime
- Uniqueness of data rows
-

FUTURE DEVELOPMENT

TODO: add user-definable options for scope of download

TODO: if title starts with brand name, remove it from title.

TODO: Sort companies by number items when scraping (shorter total time to scrape)

TODO: add categories table and brands table

TODO: Update products MSRP/TITLE/etc with new changes (provided same link)

TODO: Improve logging and usage of logging library

TODO: Add proper exception handling for concurrent futures library

TODO: Currently only takes lowest price. Add ability to pull price by color/size

USEFUL INSIGHTS

- Always comment !
Explain WHY you wrote the code, not what it does
- Methods should be ~10 lines or less
- Use shortcuts (ctrl + B, cmd + [/], ...)
- Write good asserts - error report!



main

Items for sale right now

items for sale

5,496
items for sale

3 days ago

Items added today

Items added today

6

Items added today

2 days ago

Items removed today

Items removed today

15

Items removed today

2 days ago

New items today

New items for sale today

title

brand

current_price

sizes

colors

anit

Anita Air Control Padded Sports Bra (Women's)	Anita	\$79.00	30,32,34,36,38,40,42	White,Anthracite,Pink
Anita Air Control Padded Sports Bra (Women's)	Anita	\$79.00	30,32,34,36,38,40,42	White,Anthracite,Pink
Anita Dynamix Star Control Sports Bra (Women's)	Anita	\$76.00	32,34,36,38,40	Peacock,White
Anita Dynamix Star Control Sports Bra (Women's)	Anita	\$76.00	32,34,36,38,40	Peacock,White

2 days ago

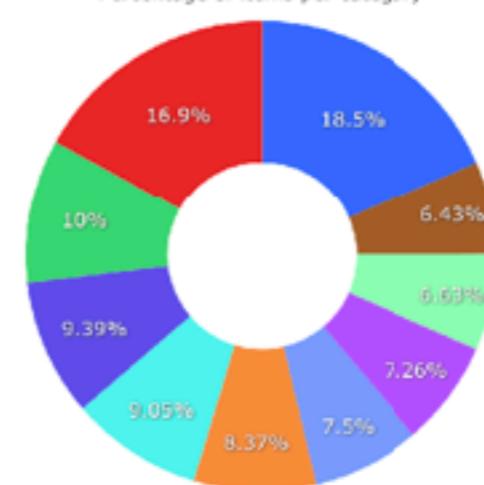
percentage products by category

Products by category

Number of categories

10

Percentage of items per category



2 days ago

Most discounted items

Highest discounts

- Ski Jackets
- Women's Jackets
- Women's Pants
- Ski Pants
- Women's Boots
- Women's Sunglasses
- Ski Mid-Layers
- Men's Boots
- Ski Hats
- Women's Hats

0 2 days ago

Percent discount per category

Products by category

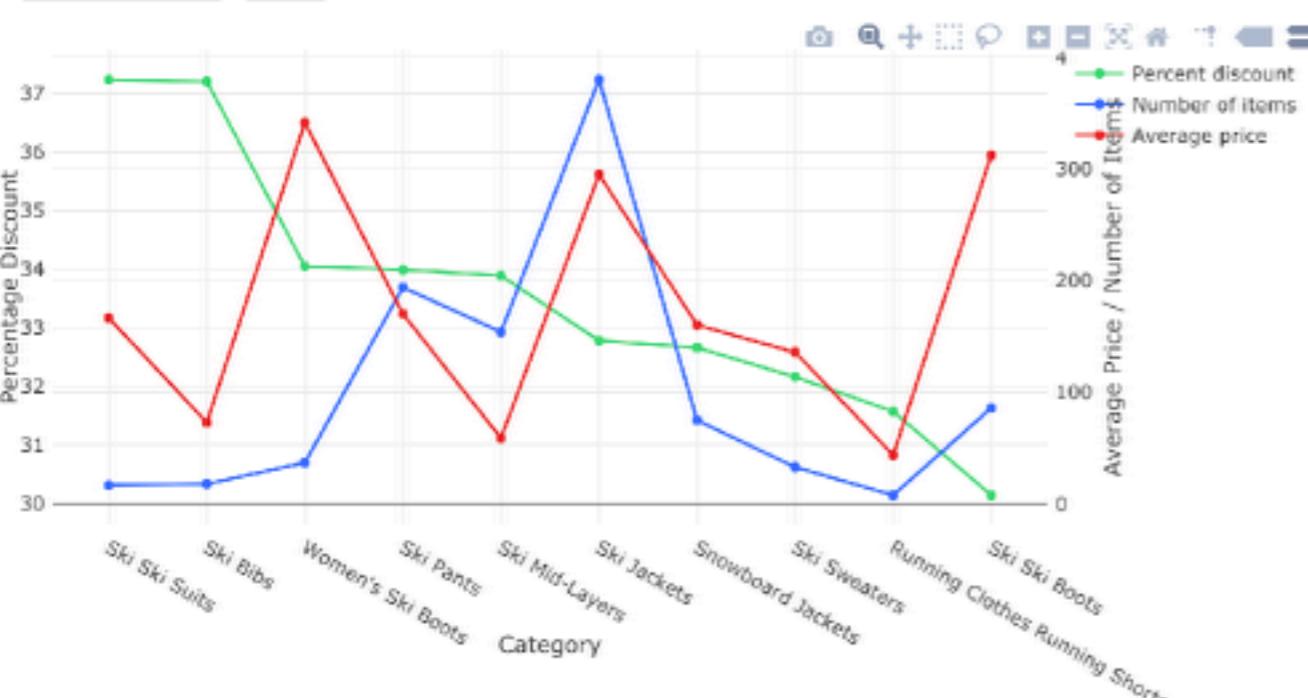
Sort by

Sort order Number OF categories

Percent discount ↓

desc ↓

10



0 2 days ago

Most discounted items

Highest discounts

discount

70 ↓

Percent Discount

Title

Link

Search...

Percent Discount	Title	Link
76.25	Obermeyer Sassy Knit Hat (Little Kids')	https://www.peterglenn.com/product/obermeyer-sassy-knit-hat-little-kids
75.02	Alpinestars Arubix Boardshorts (Men's)	https://www.peterglenn.com/product/alpinestars-arubix-boardshorts-mens
74.56	Marker Gillett Bib (Toddlers')	https://www.peterglenn.com/product/marker-gillett-bib-toddlers
74.28	Bugatchi 3/4 Coat (Men's)	https://www.peterglenn.com/product/bugatchi-34-coat-mens-0
73.34	Liquid Infinity Shell Snowboard Pant (Men's)	https://www.peterglenn.com/product/liquid-infinity-shell-snowboard-pant-mens
72.47	Montanaco Faux Lamb Ruffled Vest (Women's)	https://www.peterglenn.com/product/montanaco-faux-lamb-ruffled-vest-womens
71.42	Burton Base Camp Hybrid Short (Men's)	https://www.peterglenn.com/product/burton-base-camp-hybrid-short-mens
71.40	Mountain Force Rider III Insulated Ski Jacket (Women's)	https://www.peterglenn.com/product/mountain-force-rider-iii-insulated-ski-jacket-womens
71.32	Obermeyer Kismet Ski Jacket (Little Girls')	https://www.peterglenn.com/product/obermeyer-kismet-ski-jacket-little-girls

0 2 days ago

sankey visualization



All items for sale right now

Items for sale today



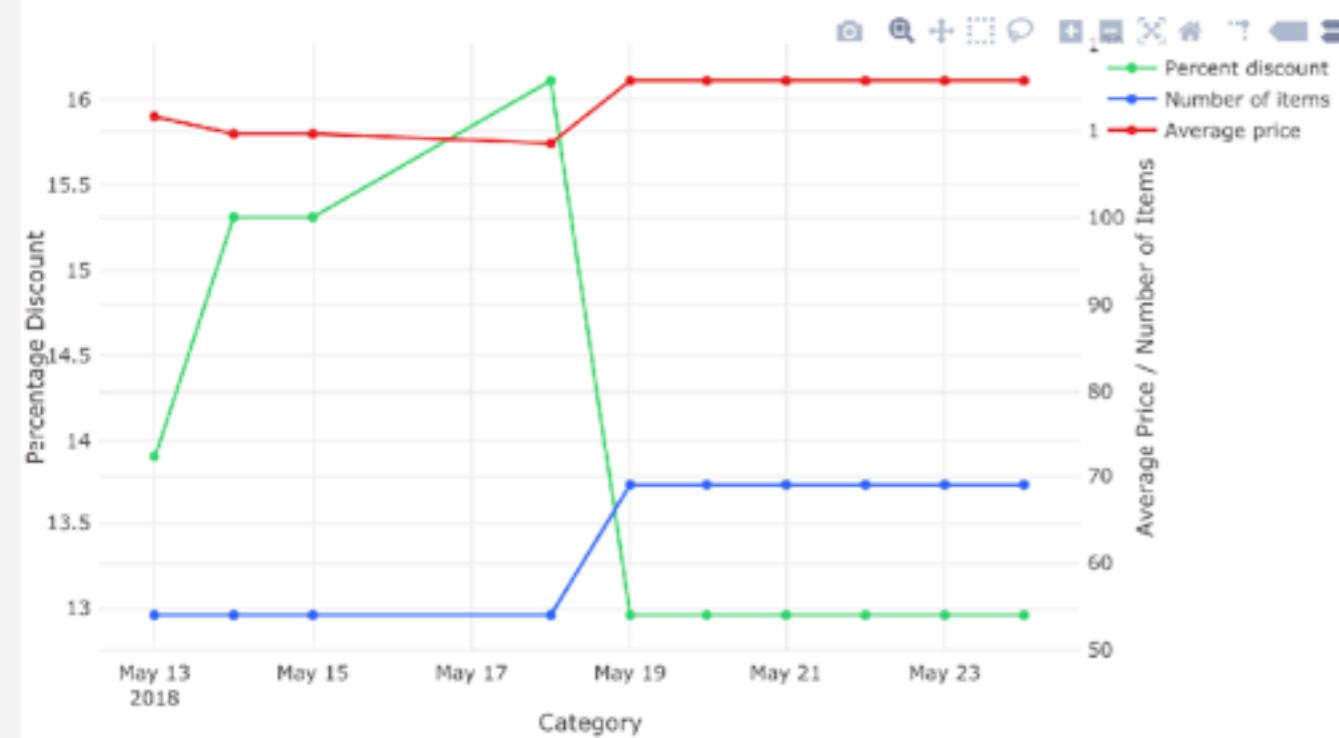
category	brand	title	current_price	msrp	sizes	colors
north face						
Men's Jackets	Vans	Vans X The North Face Torrey Coach Jacket	\$124.99	\$179.00	S,M,L,XL	TNF Black,TNF Red TNF Black,TNF Yellow TNF Black
Kids' Boots	The North Face	The North Face Shellista Lace III Boot (Kids')	\$64.99	\$75.00	10,12,13	Tagumi Brown Barolo Red
Ski Pants	The North Face	The North Face Freedom Insulated Ski Pant (Girls')				
Kids' Boots	The North Face	The North Face Shellista Extreme Boot (Girls')	\$59.99	\$69.95	1,3	Dachsund Brown Moonlight Ivory,Dark Gull Grey Cerise Pink
Women's Baselayers	The North Face	The North Face Warm Long Sleeve Crew Neck (Women's)	\$39.99	\$50.00	M	TNF Black
Ski Pants	The North Face	The North Face Snowquest Triclimate Ski Pant (Boys')	\$89.99	\$119.95	M	Cosmic Blue
Ski Pants	The North Face	The North Face Mossbud Freedom Insulated Ski Pant (Girls')	\$89.99	\$109.95	S,XL	TNF Black
Ski Mid-Layers	The North Face	The North Face Sherparazo Fleece Jacket (Boys')	\$69.99	\$84.95	XL	TNF Black
Kids' Boots	The North Face	The North Face Thermoball Shellista Boot (Girls')	\$69.99	\$80.00	1,2,3,4,6,7	TNF Black
Kids' Boots	The North Face	The North Face Thermoball Shellista Boot (Little Girls')	\$69.99	\$80.00	11,12,13	TNF Black
Kids' Jackets	The North Face	The North Face Glacier Track Fleece Jacket (Girls')	\$49.99	\$54.95	S,M,L,XL	TNF Medium Grey Heather Collar Blue,Algiers Blue
Travel And Trail Backpacks	The North Face	The North Face Pinyon Backpack (Women's)	\$59.99	\$78.95	O/S	Magic Magenta
Ski Ski Suits	The North Face	The North Face Glacier One Piece Ski Suit (Little Kids')	\$44.99	\$49.95	12M,18M,24M,6M	Sky Blue Classic Camo Print,TNF Medium Grey Heather lilac Sachet
Ski Ski Suits	The North Face	The North Face Tailout One Piece Ski Suit (Little Kids')	\$79.99	\$98.95	0_3MO,3_6MO,6_12MO	Bright Cobalt Blue
Ski Ski Suits	The North Face	The North Face Tailout One Piece Ski Suit (Little Kids')	\$79.99	\$98.95	0_3MO,3_6MO	Bright Cobalt Blue

Category discount by date

Products by category

Category

Women's Running Shoes



2 days ago

Box plot distribution of price and MSRP right now

