



JOHNS HOPKINS  
CAREY BUSINESS SCHOOL

# Lecture 5

## **BU.330.740 Large Scale Computing on the Cloud**

Minghong Xu, PhD.  
Associate Professor

# Reflections



## » Recommendation Systems

- Search vs Discovery
  - Long tail, power law, 80-20 rule → *silent majority.*
  - Main task: predict attitudes
  - Two entities: users and items
  - Two basic approaches: content-based vs collaborative-filtering-based
  - Models:
    - KNN
    - K-means
    - Item-item
    - Factorization machine
- network externalities.*



# Factorization Machine

- » Decompose the prediction of user-item attitudes into a few factors
- » Idea similar to regression models, factor models

*Fama and French 5-factor model.  
↳ read about this*



# Model Deployment

## » Endpoint

- Refer to a deployed model or service that can handle requests
- In AWS SageMaker, an endpoint allows you to make real-time predictions using a trained ML model
- Another example: DeepSeek

## » Remember to terminate all endpoints!



# Learner Account Limitations

- » Learner account is not stable for model deployment
  - 6/180 accounts denied
- » Check all limitations under:  
Launch AWS Academy Learner Lab->Service usage and other restrictions
- » Register account directly with AWS for full, unlimited services



# Today's Agenda

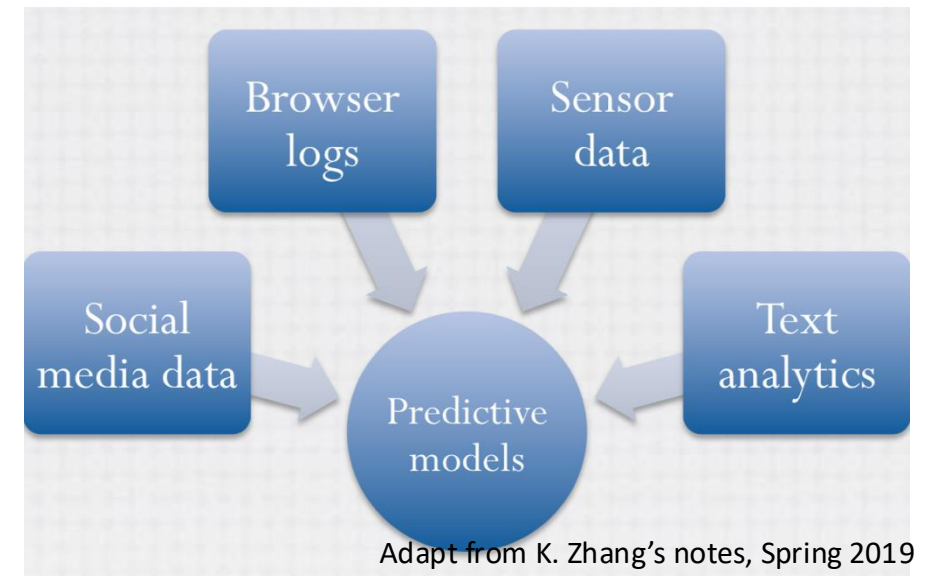
- » Mining of massive dataset
- » Computer vision and business applications
- » Lab 4: AWS Rekognition
  - Custom labels
  - Object detection
  - Facial analysis
  - Text in image
- » No model deployment today



# Mining of Massive Dataset

# Case: Marketing

- » Big data is used to better understand customers and their behaviors and preferences
  - Target: very accurately predict when one of their customers will expect a baby
  - Wal-Mart can predict what products will sell
  - Car insurance companies understand how well their customers actually drive





# Applications and Data



## » Applications

- Sentiment analysis
- Recommender system

## » Data

- Search and user logs
- Customer transaction records
- User generated content

## » Data Characteristics

- Structured and unstructured containing text and image

Web Logs

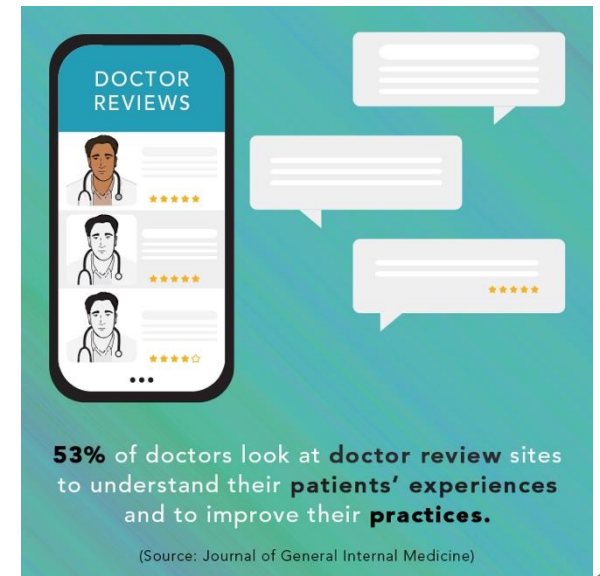
Time	Computer	Group	User	Caption	URL
2012-01-06 14:24:57	Kevin	Admin	Kevin	eTinySoft Total Video Converter, Total V...	www.effectmatrix.com/index.htm
2012-01-06 14:22:25	Kevin	Admin	Kevin	PowerPoint convert software - E.M. Pow...	www.effectmatrix.com/PowerPoint-Video-Converter/in...
2012-01-06 14:22:03	Kevin	Admin	Kevin	eTinySoft Total Video Converter, Total V...	www.effectmatrix.com/index.htm
2012-01-06 14:20:41	Kevin	Admin	Kevin	Free PowerPoint to AVI - Free PowerPoi...	www.effectmatrix.com/PowerPoint-Video-Converter/Fr...
2012-01-06 14:20:02	Kevin	Admin	Kevin	ppt to video converter - Google Search ...	https://www.google.com/#client=psy-ab&hl=en&new...
2012-01-06 14:19:08	Kevin	Admin	Kevin	PowerPoint to Video Converter, Free do...	www.xilisoft.com/ppt-software.html
2012-01-06 14:18:46	Kevin	Admin	Kevin	Xilisoft Video Converter, DVD Ripper, iPa...	www.xilisoft.com
2012-01-06 14:17:28	Kevin	Admin	Kevin	PowerPoint to Video, convert PowerPoi...	www.dvd-ppt-slideshow.com/ppt-to-video/
2012-01-06 14:17:17	Kevin	Admin	Kevin	ppt to video converter - Google Search ...	https://www.google.com/#client=psy-ab&hl=en&new...
2012-01-06 14:15:39	Kevin	Admin	Kevin	ppt to video converter - Google Search ...	https://www.google.com/#client=psy-ab&hl=en&new...
2012-01-06 14:14:44	Kevin	Admin	Kevin	PPT to DVD Converter: Convert PowerP...	www.wondershare.com/pro/ppt2dvd-pro.html
2012-01-06 14:03:27	Kevin	Admin	Kevin	4Videosoft Video Converter Tools, Best ...	www.4videosoft.com
2012-01-06 11:54:32	Kevin	Admin	Kevin	Download Video Tools, DVD Tools, Pow...	www.leawo.com
2012-01-06 11:52:44	Kevin	Admin	Kevin	Freemove Freeware - Free video convert...	www.freemove.com
2012-01-06 11:50:32	Kevin	Admin	Kevin	Freemove Freeware - Free video convert...	www.freemove.com
2012-01-06 11:18:27	Kevin	Admin	Kevin	Cucusoft Video Converter, DVD to ipod, L...	www.google.com
2012-01-06 11:08:11	Kevin	Admin	Kevin		
2012-01-06 10:44:40	Kevin	Admin	Kevin		
2012-01-06 10:41:12	Kevin	Admin	Kevin		
2012-01-06 10:23:12	Kevin	Admin	Kevin		

	A	B	C	D	E	F
1	Table: Transaction				±: Credit, -: Debit	
2	Transaction ID	Transaction Date Time	User ID	Account ID	Amount	Account Balance
3	10000001	01/04/2012 09:10:19	2	1	3100.00	4,300.21
4	10000002	01/04/2012 11:10:19	4	3	5800.00	6,412.44
5	10000003	01/04/2012 12:10:19	3	4	1200.00	307.85
6	10000004	01/04/2012 13:10:19	1	5	2500.00	229.87
7	10000005	02/04/2012 09:10:19	5	1	-50.00	4,250.21
8	10000006	02/04/2012 11:10:19	3	3	-100.00	612.44
9	10000007	02/04/2012 14:10:19	1	6	810.00	-99,119.91
10	10000008	03/04/2012 09:10:19	3			
11	10000009	03/04/2012 11:10:19	1			
12	10000010	03/04/2012 14:10:19	5			



# Case: Healthcare

- » Big data analytics allow us to monitor and predict the developments of epidemics and disease outbreaks
- » Streaming data from blood sugar monitoring sensors can be used to predict poorly controlled diabetes
- » Doctors' reviews can be used to analyze patient satisfaction



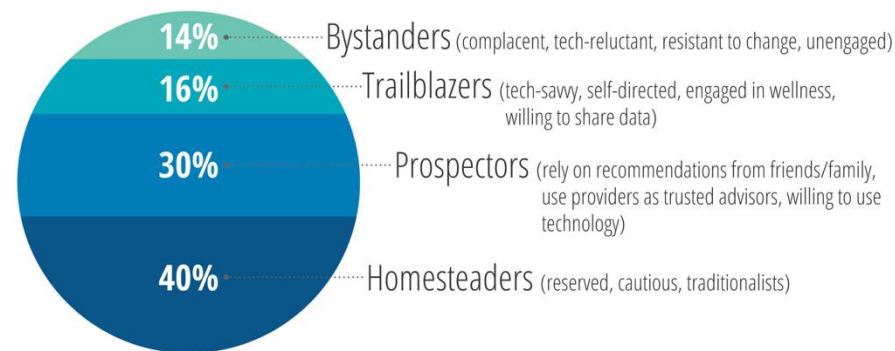
# Healthcare (Cont.)



- » Market segmentations can be used to achieve personalized care
  - Understand attitudes and behaviors to attract, retain, and engage consumers
- » Deloitte's healthcare consumer segment profiles:  
<https://www2.deloitte.com/us/en/insights/industry/healthcare/healthcare-consumer-patient-segmentation.html>

FIGURE 1

Distribution of segments in the Deloitte 2018 Survey of US Health Care Consumers



Source: Deloitte 2018 Survey of US Health Care Consumers.

Deloitte Insights | [deloitte.com/insights](https://deloitte.com/insights)

# Applications and Data



## » Applications

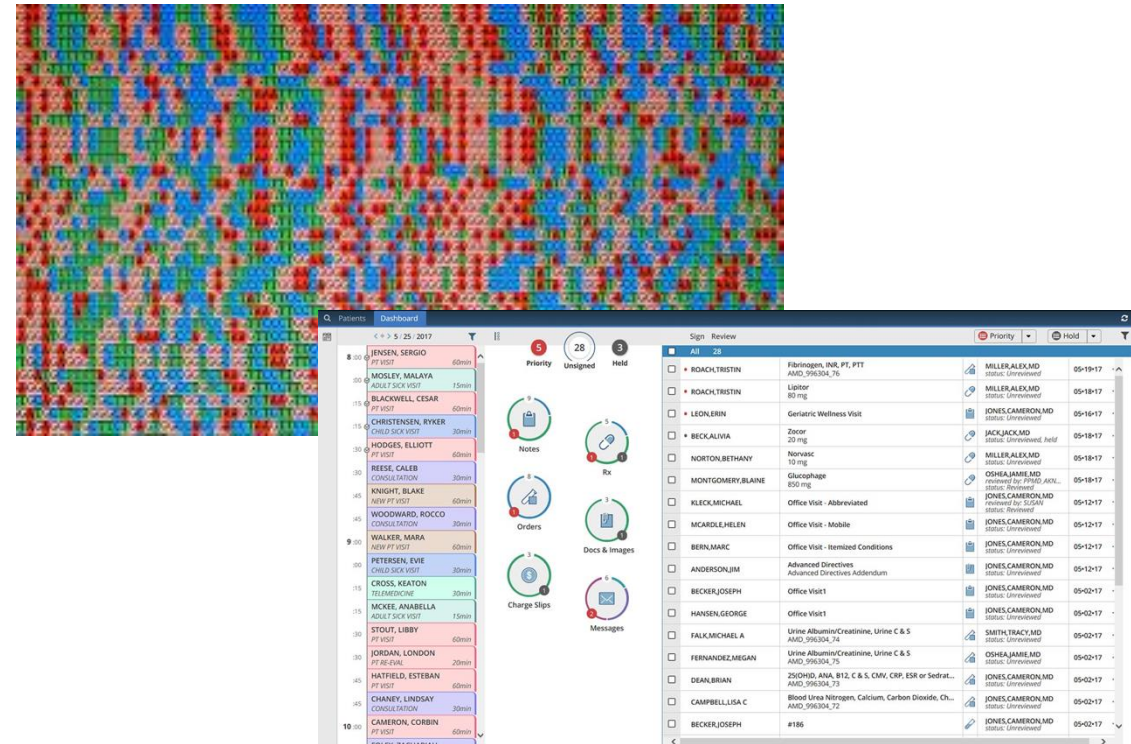
- Human and plant genomics
- Decision support system
- Patient segmentation

## » Data

- Genomics and sequence data
- Electronic medical records (EMR)
- Health and patient social media

## » Data Characteristics

- Disparate but highly linked content, person-specific content, and ethics issues



# Case: Security and Public Safety



## » Applications

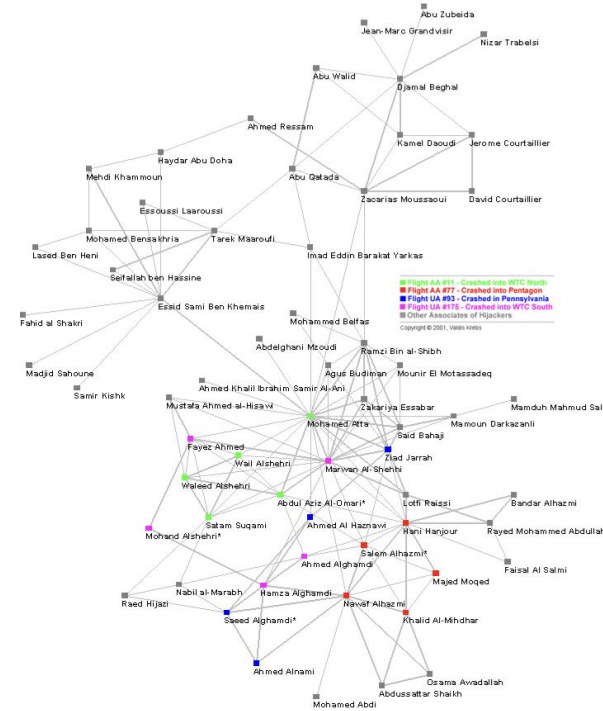
- Detect terrorists
- Track cyber attacks and botnets
- Spam filter

## » Data

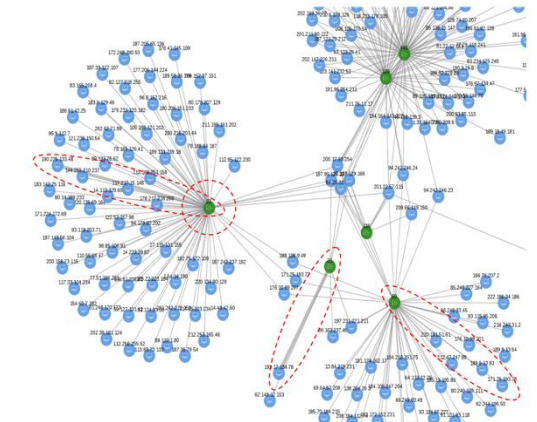
- Criminal records
- Crime maps
- Criminal networks
- Web trails

## » Data Characteristics

- Personal identity information, deceptive, multilingual



Krebs, V. E. (2002). Mapping networks of terrorist cells. *Connections*, 24(3), 43-52.







# Other Cases

## » Sports

- Use video analytics to track the performance of every player
- Use sensor technology in sports equipment to get feedback on games

## » Supply Chain

- Optimize stock based on predictions generated from social media data, web search trends, and weather forecasts
- Geographic positioning and radio frequency identification sensors are used to track goods or delivery vehicles and optimize routes by integrating live traffic data

## » Finance, e.g. fraud detection

## » Smart city, e.g. optimize energy grids



# Discussion Time

» *What data is your industry using?*

- Industry: you interned or worked
- Or data that can be potentially used

» *How is the data used?*

- Or how you propose to use it

# Computer Vision and Applications



# Overview

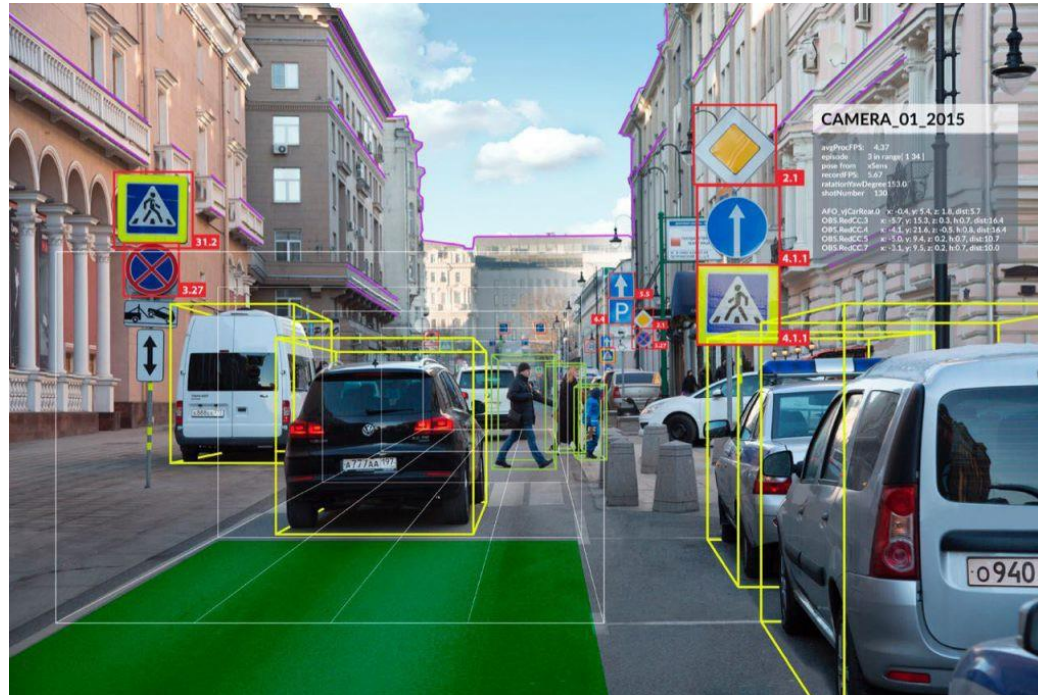


- » Computer Vision: automated extraction of information from digital images
- » Applications
  - Public safety and home security
  - Content management and analysis
  - Autonomous driving
  - Medical imaging
  - Manufacturing process control

# Image Applications I



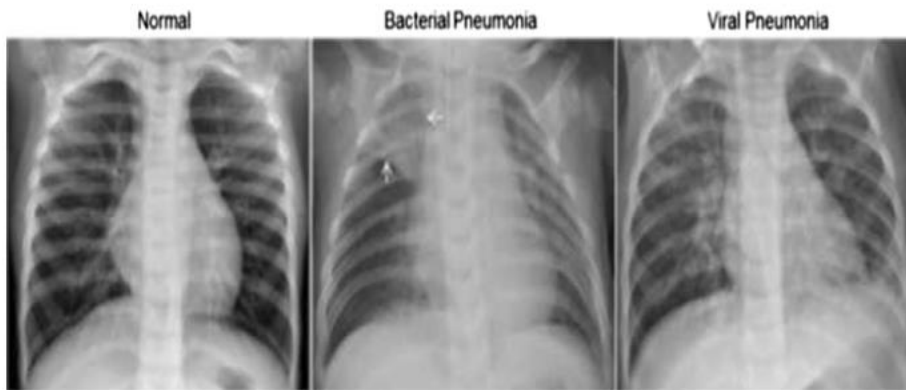
- » Tesla Autopilot
- » <https://www.linkedin.com/pulse/teslas-use-ai-revolutionary-approach-car-technology-alexander-stahl/>



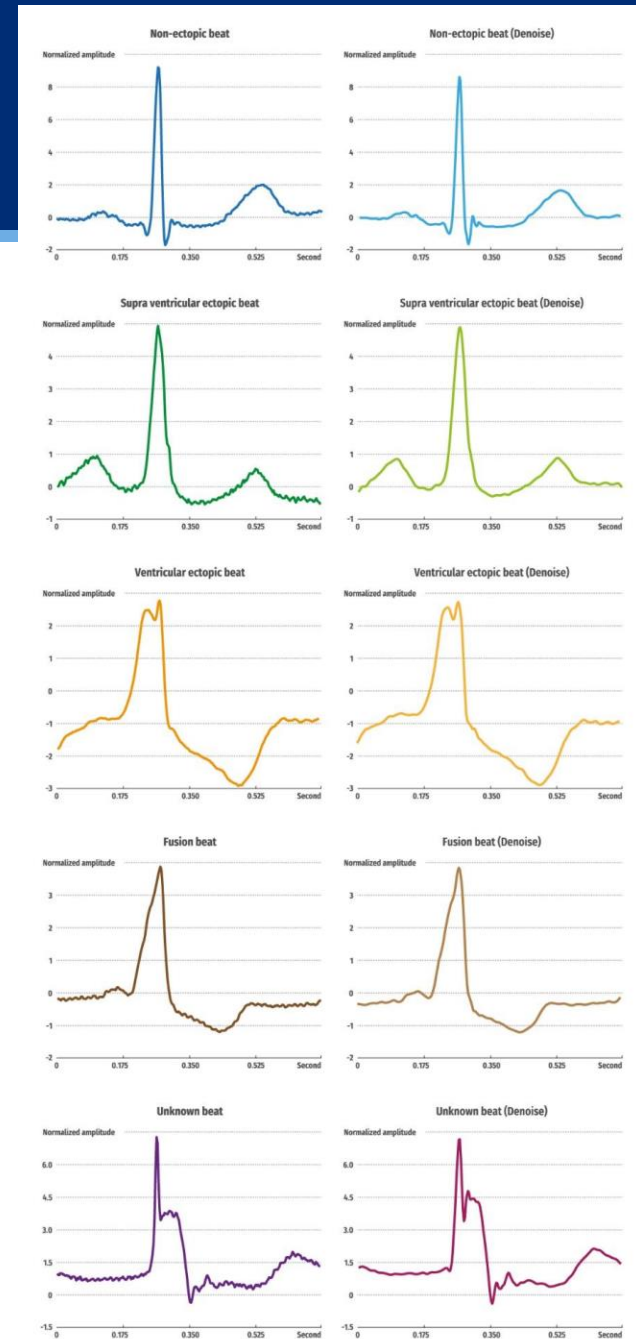
# Image Applications II

## » Medical Imaging

- Chest X-ray
- Heart ECG



**Fig. 1** Examples of chest X-rays [3]



# Image Applications III



- » Satellite image, parking lot traffic
- » <https://newsroom.haas.berkeley.edu/how-hedge-funds-use-satellite-images-to-beat-wall-street-and-main-street/>





# Image Applications IV



## » The Nature Conservancy Fisheries Monitoring

- Kaggle competition: <https://www.kaggle.com/c/the-nature-conservancy-fisheries-monitoring>

## » Google: AI camera system, computer vision to monitor health of fish populations

-  ALB: Albacore tuna (*Thunnus alalunga*)
-  BET: Bigeye tuna (*Thunnus obesus*)
-  DOL: Dolphinfish, Mahi Mahi (*Coryphaena hippurus*)
-  LAG: Opah, Moonfish (*Lampris guttatus*)
-  SHARK: Various: Silky, Shortfin Mako
-  YFT: Yellowfin tuna (*Thunnus albacares*)



Fish images are not to scale with one another



<https://www.dailymail.co.uk/sciencetech/article-8066017/Alphabet-unveils-AI-camera-monitors-fish-populations-goal-feeding-humanity.html>

# Image Applications V



» Image style transfer

» <https://github.com/lengstrom/fast-style-transfer>

**Style**  
*The Starry Night*,  
Vincent van Gogh,  
1889



**Style**  
*The Muse*,  
Pablo Picasso,  
1935





# Image Applications VI



- » Image “sentiment” analysis
- » Visual listening in: Extracting brand image portrayed on social media (Liu, Liu, et al. 2020 *Marketing Science* 39.4: 669-686)
  - Mine visual content posted on Instagram on brand attributes: glamorous, rugged, healthy, fun
  - Build training data from Flickr by query attribute word and collect top 200



Figure 1 Sample images from Instagram hashtagged with brands



Figure 3 Sample images of different brand attributes and their antonyms in our data sets

# Image Applications VII

» Image-based recommendations on styles and substitutes (McAuley, Julian, et al. 2015 *Proceedings of the 38th international ACM SIGIR conference on research and development in information retrieval*)

- Model visual relationships: alternatives, complementary
- Recommending which clothes and accessories will go well together
- Amazon dataset of 180 million relationships between almost 6 million objects



Figure 8: Outfits generated by our algorithm (Women's outfits at left; Men's outfits at right). The first column shows a 'query' item that is randomly selected from the product catalogue. The right three columns match the query item with a top, pants, shoes, and an accessory, (minus whichever category contains the query item).





# Discussion Time

- » *What image data is your industry using?*
- » *How is the data used?*
  - Or how you propose to use it

# Tasks: Image Analysis



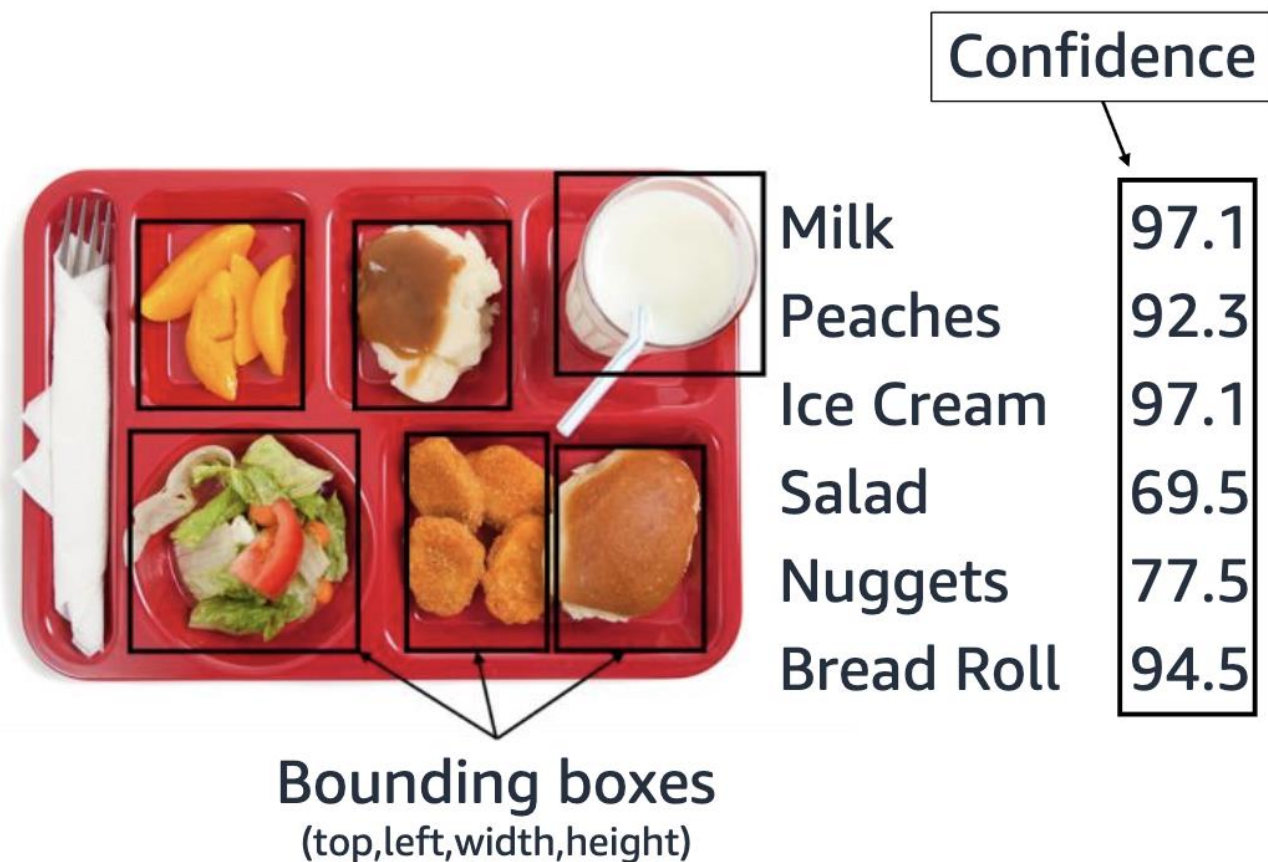
- » Object classification: identify and categorize objects within an image into predefined classes



Food?  
Breakfast?  
Lunch?  
Dinner?

# Image Analysis (Cont.)

- » Object detection: identify and locate objects within an image or video



# Image Analysis (Cont.)



- » Object segmentation/semantic segmentation: identify and segment objects in an image at the pixel level



Milk  
Peaches  
Ice Cream  
Salad  
Nuggets  
Bread Roll

# Comparison



Feature	Object Classification	Object Detection	Object Segmentation
Definition	Identifies what object is present in an image	Detects objects and their locations with bounding boxes	Identifies and outlines objects at the pixel level
Output	Single label per image (e.g., "cat")	Bounding boxes around detected objects	Pixel-wise masks showing object boundaries
Precision	Least precise (only recognizes presence)	Moderate precision (box-based localization)	Most precise (detailed object shape)
Use Cases	Image categorization (e.g., "dog" vs. "cat")	Face detection, pedestrian detection	Medical image analysis, autonomous driving
Example	"This is a car."	"This is a car located at (x, y, w, h)."	"This is a car, and these are its exact pixel boundaries."
Complexity	Simple (single classification per image)	Moderate (detects multiple objects in an image)	High (detailed segmentation of objects)

Imp.



# Tasks: Video Analysis

## » Instance tracking

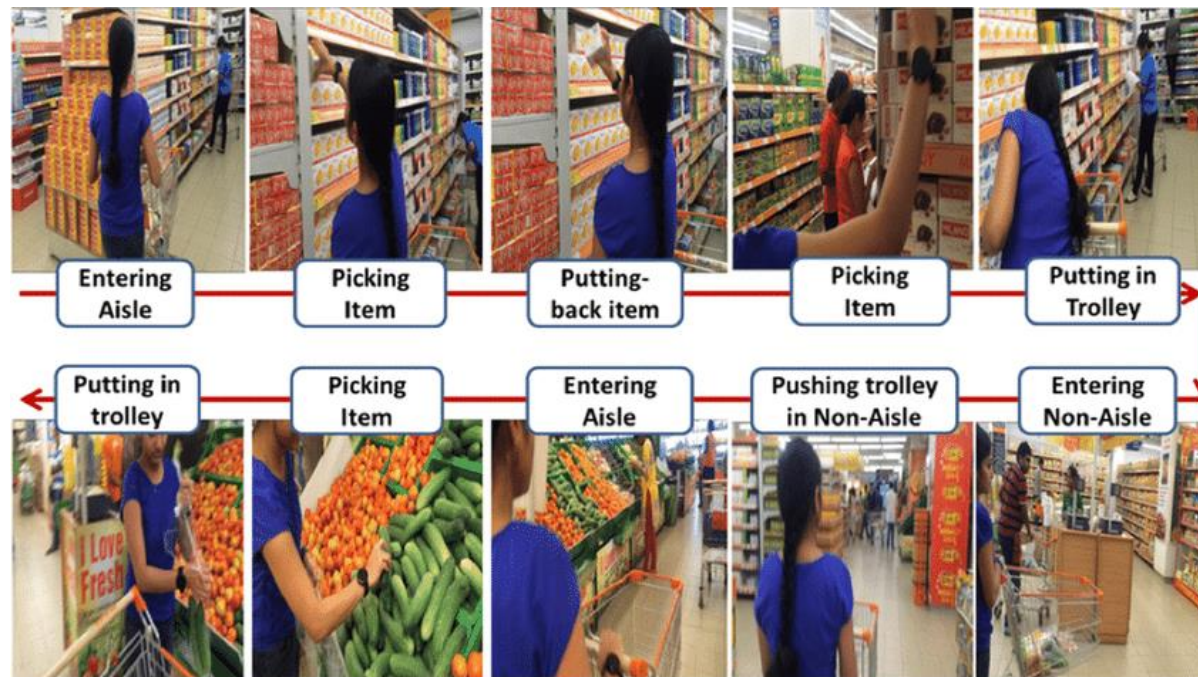
- Capture the path of people in the scene
- For example, use the movement of athletes during the game for post-game analysis



# Video Analysis (Cont.)

## » Action recognition

- Analyze shopper behavior and density in your retail store by studying the path that each person follows

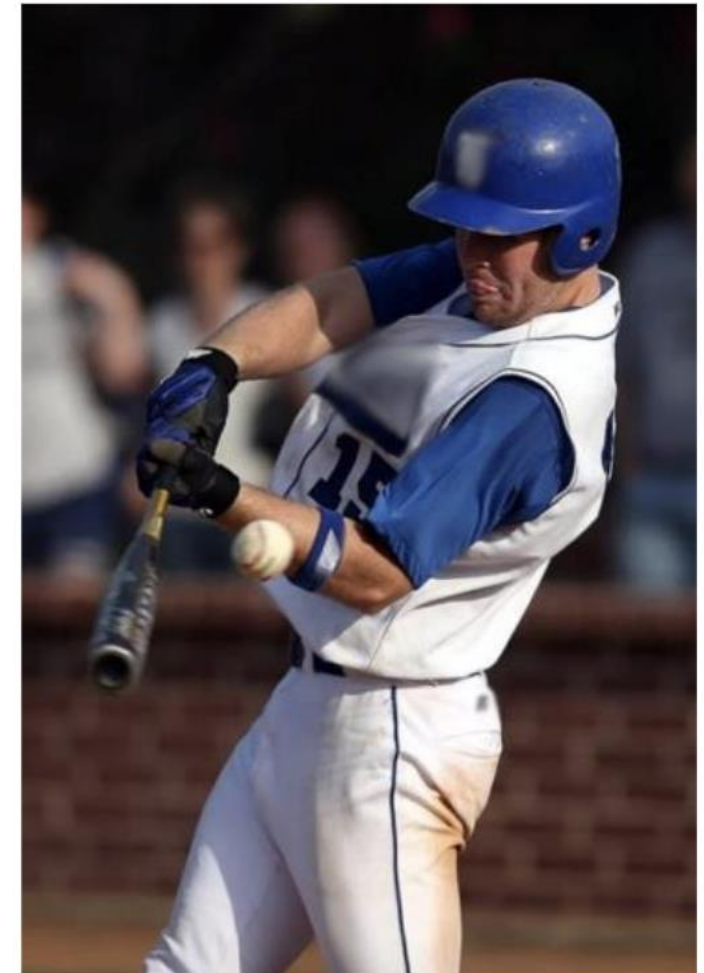


# Video Analysis (Cont.)



## » Motion estimation

- Batter's accuracy, pitcher's pitching style, type of pitch, inning, batter's performance vs specific pitcher, etc.
- For example, speed of ball leaving the bat and trajectory -> probability of a home run, possible foul ball into the crowd





# Lab 4



- » AWS Rekognition
- » Custom labels take a long time to train
  - 20-30 minutes
  - Cost ~\$2
- » No extension today
  - Please start your project

# Next Week



- » Natural language processing on cloud
- » AI-assisted programming

# References



» AWS Academy Machine Learning Foundations