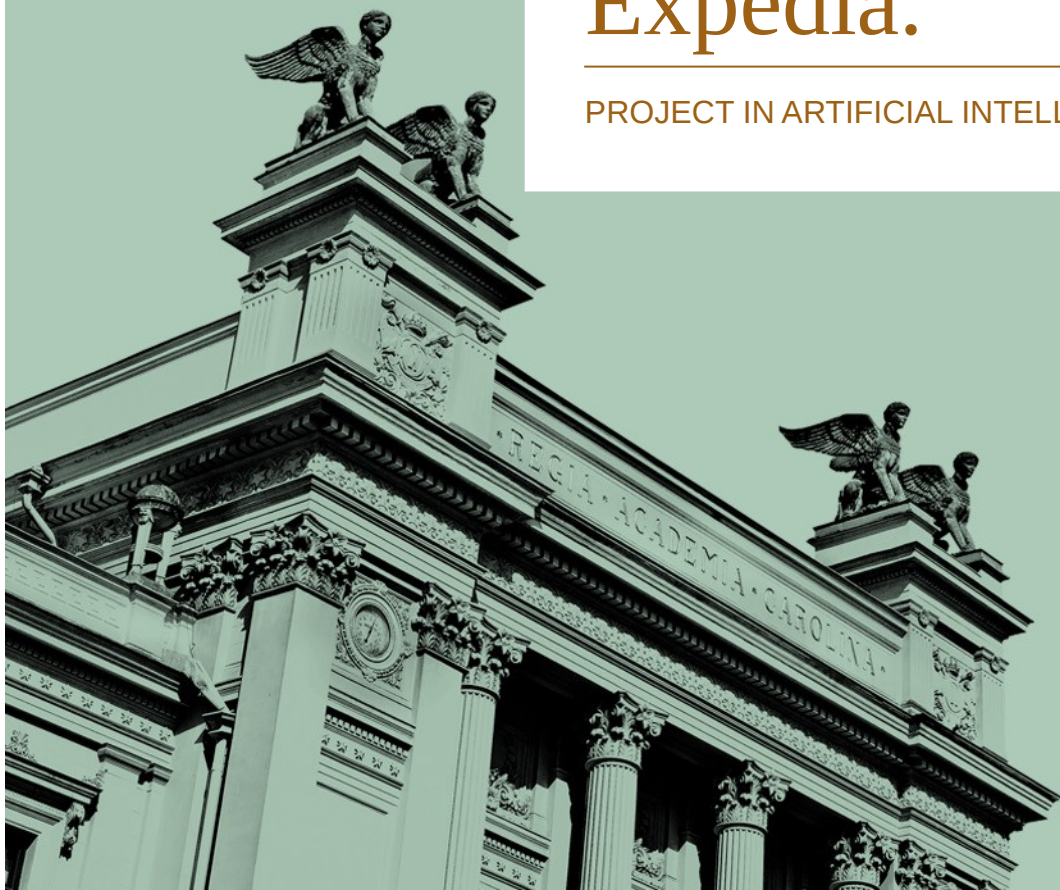




LUND  
UNIVERSITY

# Classification of hotels for Expedia.

PROJECT IN ARTIFICIAL INTELLIGENCE - EDAN70



# Introduction

---

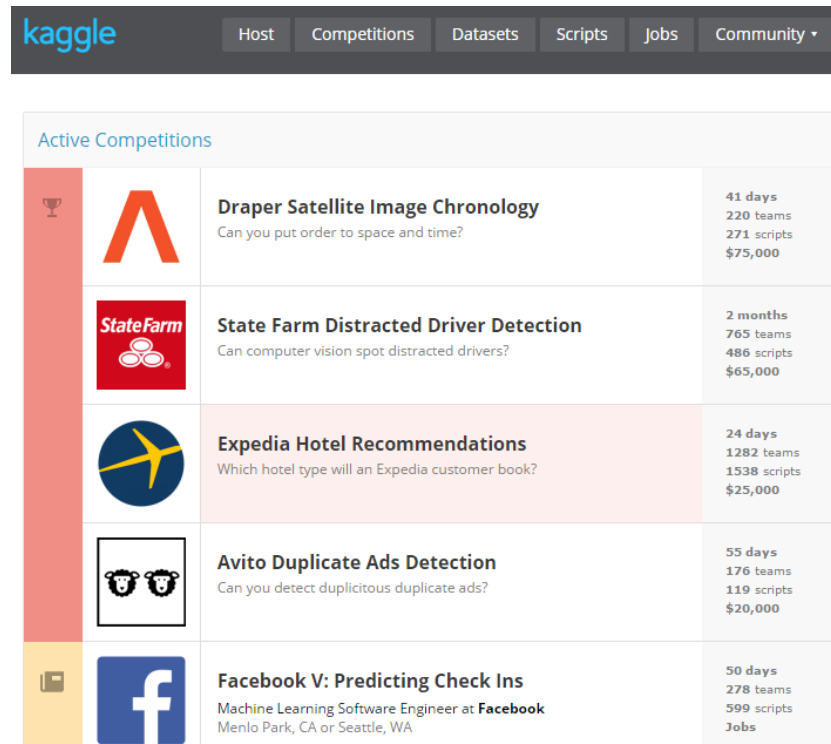
- Who we are.
- Kaggle.com
- Our main problem. Expedia
- Random Forest Classifier.
- Expedia and workflow.
- Conclusions.








# Kaggle.com

---

- Users from all over the world compete to produce the best machine learning models.
- Submissions, Scripts, Leaderboards.



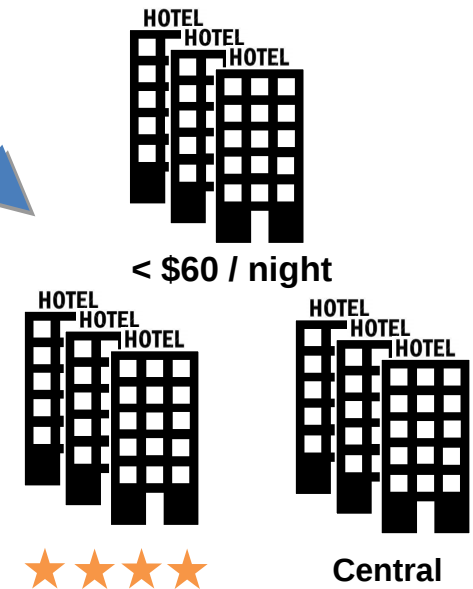
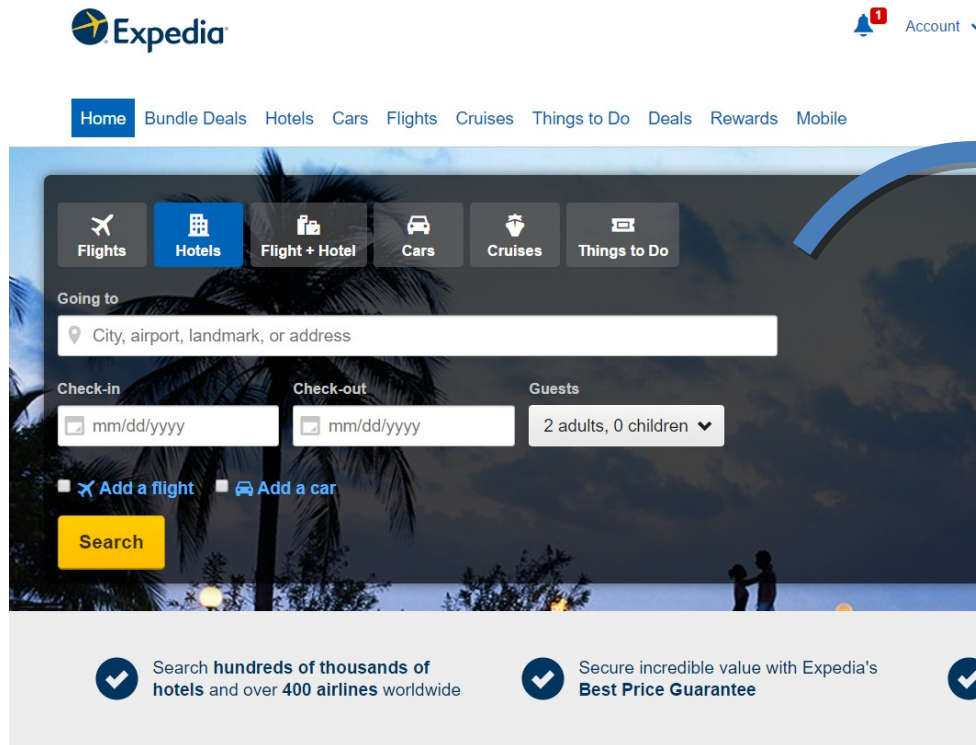
The screenshot shows the Kaggle website's navigation bar and a list of active competitions. The navigation bar includes the Kaggle logo and links for Host, Competitions, Datasets, Scripts, Jobs, and Community. The active competitions section lists five competitions with their respective logos, titles, descriptions, and statistics.

Competition	Duration	Teams	Scripts	Prize
 <b>Draper Satellite Image Chronology</b> Can you put order to space and time?	41 days	220 teams	271 scripts	\$75,000
 <b>State Farm Distracted Driver Detection</b> Can computer vision spot distracted drivers?	2 months	765 teams	486 scripts	\$65,000
 <b>Expedia Hotel Recommendations</b> Which hotel type will an Expedia customer book?	24 days	1282 teams	1538 scripts	\$25,000
 <b>Avito Duplicate Ads Detection</b> Can you detect duplicitous duplicate ads?	55 days	176 teams	119 scripts	\$20,000
 <b>Facebook V: Predicting Check Ins</b> Machine Learning Software Engineer at <b>Facebook</b> Menlo Park, CA or Seattle, WA	50 days	278 teams	599 scripts	Jobs



# Expedia

- The problem – Expedia.



# Ho Chi Minh City: 530 hotels



Sort By: Price Guest Rating Vacation Rentals Hotel Name Hotel Class Recommended More

**45%** booked  
Ho Chi Minh City is a popular location on your dates.  
Hint: you can get a lower price on other dates  
[Try searching one week later](#)

Search by hotel name

Filter hotels by

**NEW!** Free Cancellation

Hotel Class

- ★★★★★ 5 Stars (18)
- ★★★★ 4 Stars (52)
- ★★★ 3 Stars (215)
- ★★ 2 Stars (222)
- ★ 1 Star (23)

Price Per Night

- Less than \$75 (95)
- \$75 to \$124 (20)
- \$125 to \$199 (5)
- \$200 to \$299 (2)
- Greater than \$300 (1)

Search Nearby

- Opera House
- War Remnants Museum
- Saigon Square
- Diamond Plaza
- Reunification Palace

[Show more](#)

Neighborhood

- Ho Chi Minh City (and vicinity)
- Ben Thanh Market
- Bui Thi Xuan
- Consulate Area
- District 3

Find what you want faster, filter by price or neighborhood.  
[Got it.](#)

Hotel avg \$81    3 star avg \$69    4 star avg \$131    5 star avg \$307

**Today 24% off**



**Saigon Odyssey 1 Hotel** ★★  
Pham Ngu Lao [Map](#)  
1-888-284-5744 • Expedia Rate  Free Cancellation

**Good! 3.6/5**  
(28 reviews)  
We have 5 left at  
~~\$46~~ **\$35** avg/night  
**Sale!**  
Hurry! Offer ends in 08:12:47  
Earn 2,248 points

**Today 25% off**



**Saigon Domaine Luxury Residences** ★★★★★  
#5 Guest Rated  
Ho Chi Minh City [Map](#)  
1-888-287-9053 • Expedia Rate

**Wonderful! 4.5/5**  
(152 reviews)  
~~\$198~~ **\$149** avg/night  
**Sale!**  
Earn 9,620 points

**Get an extra 10% or more off select hotels with Member Pricing!**

Already a member? [Sign In](#)

You agree to receive deals and offers from Expedia, and may unsubscribe at any time

**Today 31% off**



**Lan Lan Hotel 2** ★★★  
#75 Guest Rated  
Ben Thanh Market [Map](#)  
1-888-272-4858 • Expedia Rate  Free Cancellation

**Good! 3.8/5**  
(166 reviews)  
In high demand!  
We have 2 left at  
~~\$75~~ **\$52** avg/night  
**Sale!**  
Earn 6,130 points



ND  
RSITY

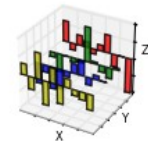
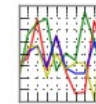
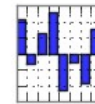
# Tools

---

- Python 64-bit
- A lot of RAM
- Pandas – parsing data into data structures
- NumPy – scientific computing package
- Scikit learn – Machine Learning library, built on SciPy, NumPy and matplotlib

pandas

$$y_{it} = \beta' x_{it} + \mu_i + \epsilon_{it}$$



LUND  
UNIVERSITY

# Expedia - dataset

---

- 24 column in training
- 22 column in testing (no is\_booking, no hotel\_cluster)
- Most of the columns are integers or floats
- Output hotel cluster ID integer range from 1-99

	date_time	site_name	posa_continent	user_location_country	user_location_region	user_location_city	...
0	2014-08-11 07:46:59	2	3	66	348	48862	...
1	2014-08-11 08:22:12	2	3	66	348	48862	...
2	2014-08-11 08:24:33	2	3	66	348	48862	...
3	2014-08-09 18:05:16	2	3	66	442	35390	...
4	2014-08-09 18:08:18	2	3	66	442	35390	...



# Expedia - workflow

- Understanding dataset



Going to

srch\_destination\_type\_id, hotel\_continent, hotel\_country, and hotel\_market

Check-in  Check-out  Guests 2 adults, 0 children ▼

srch\_ci srch\_co are filled with dates  
srch\_adults\_cnt, srch\_children\_cnt, and srch\_rm\_cnt is number of guests and rooms

**Add a flight** Global Sites

Add a flight maps to the is\_package field

site\_name – Expedia point of sale (Expedia.com, Expedia.se, ...)

posa\_continent – ID of continent associated with site\_name





# Expedia – Hotel Clusters

## Filter hotels by

**NEW!** Free Cancellation

## Hotel Class

- ★★★★★ 5 Stars (1)
- ★★★★ 4 Stars (5)
- ★★★ 3 Stars (64)
- ★★ 2 Stars (31)
- ★ 1 Star (1)

## Price Per Night

- Less than \$75 (14)
- \$75 to \$124 (10)
- \$125 to \$199 (0)
- \$200 to \$299 (0)
- Greater than \$300 (0)

## Neighborhood

- Ho Chi Minh City (and vicinity)
- Ben Thanh Market
- Bui Thi Xuan
- Consulate Area

## Property Type

- Hotel (100)
- Apartment (1)
- Apart-hotel (1)

## Meal Plans

- All Inclusive
- Breakfast
- Full Board
- Half Board

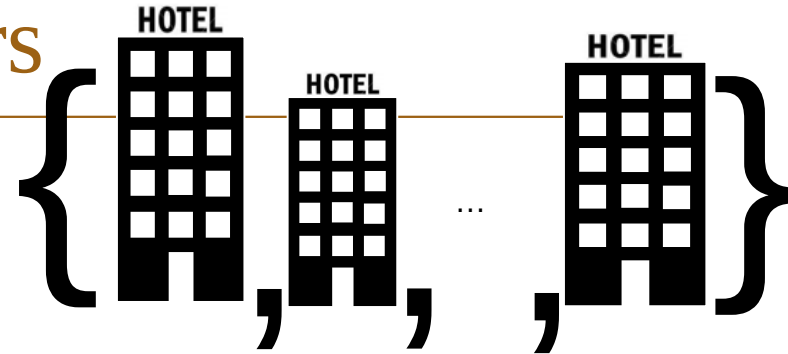
## Amenities

- High-speed Internet (102)
- Air conditioning (98)
- Swimming pool (5)
- Babysitting service (6)
- Business services (50)

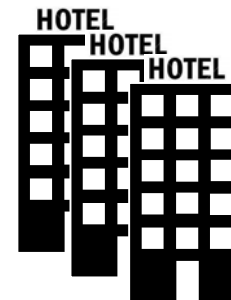
[Show more](#)

## Accessibility

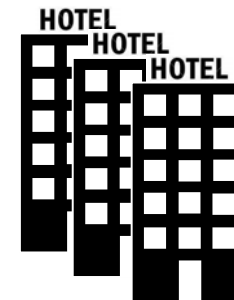
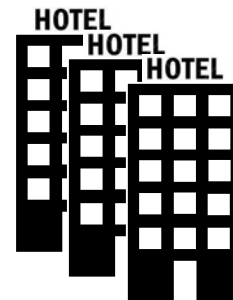
- Accessible bathroom (10)
- In-room accessibility (11)



Useful! Expedia can much quicker at an earlier stage filter the hotels



< \$60 / night



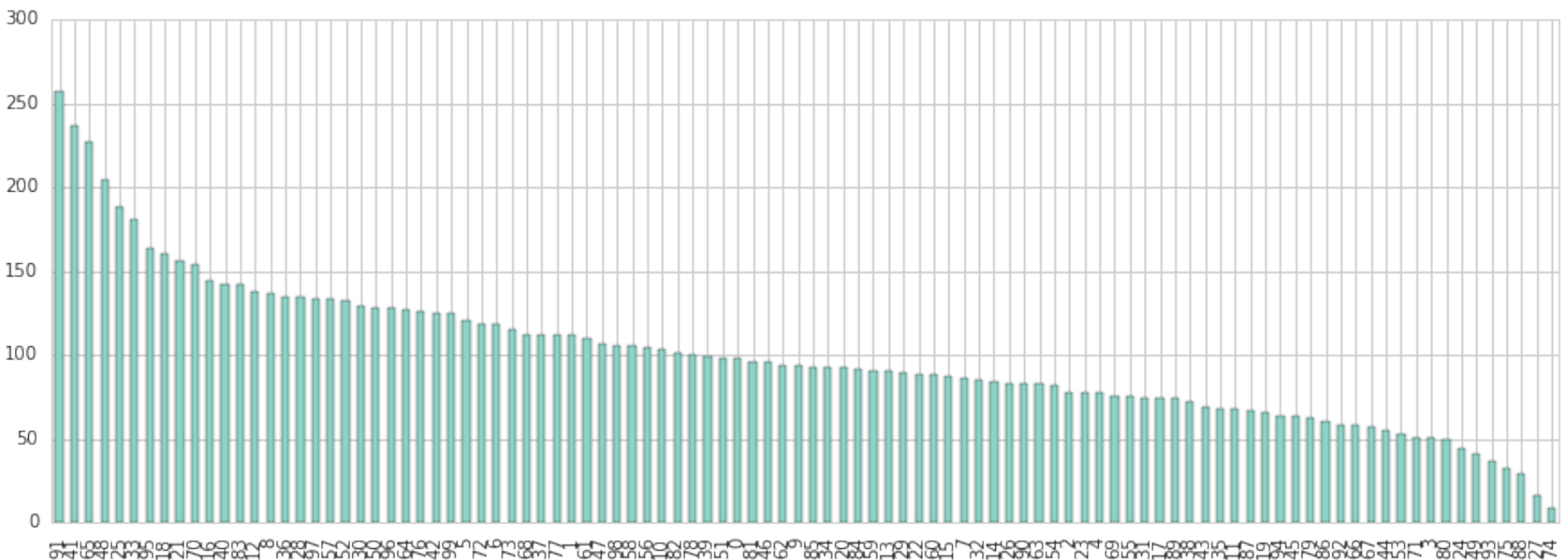
Central



LUND  
UNIVERSITY

# Expedia – most frequent hotel clusters

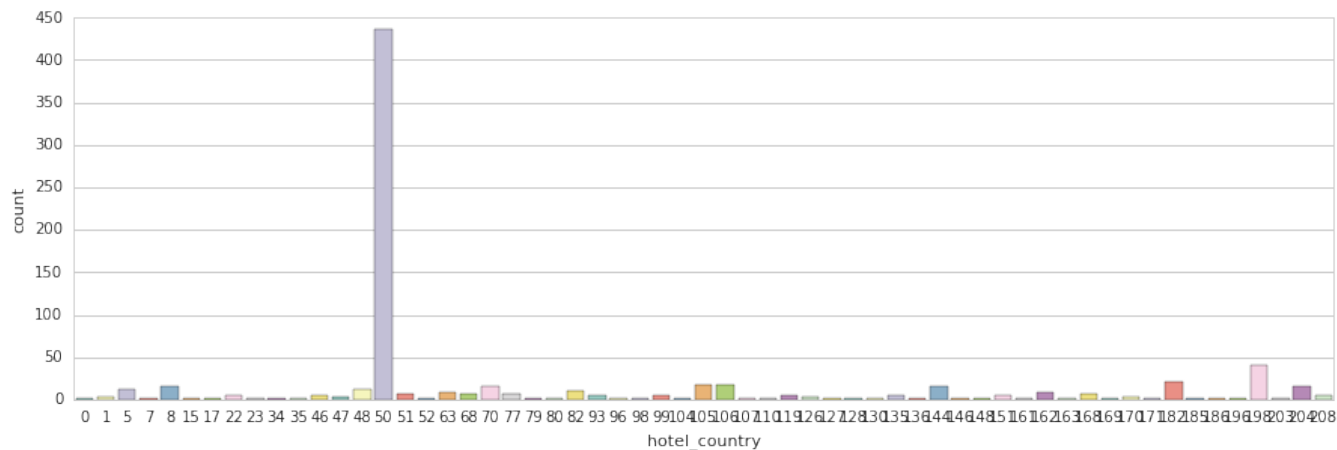
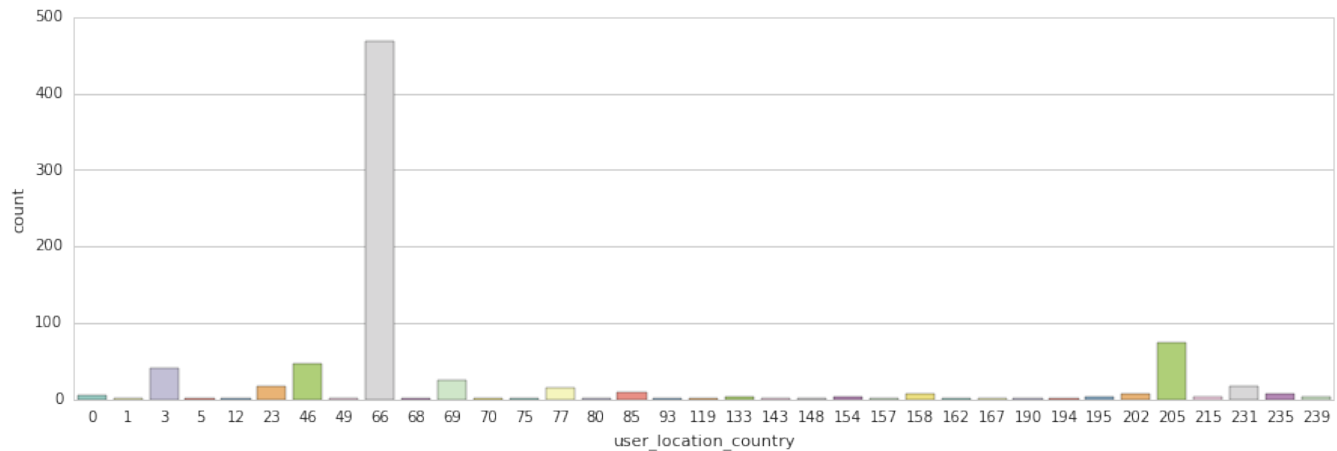
---



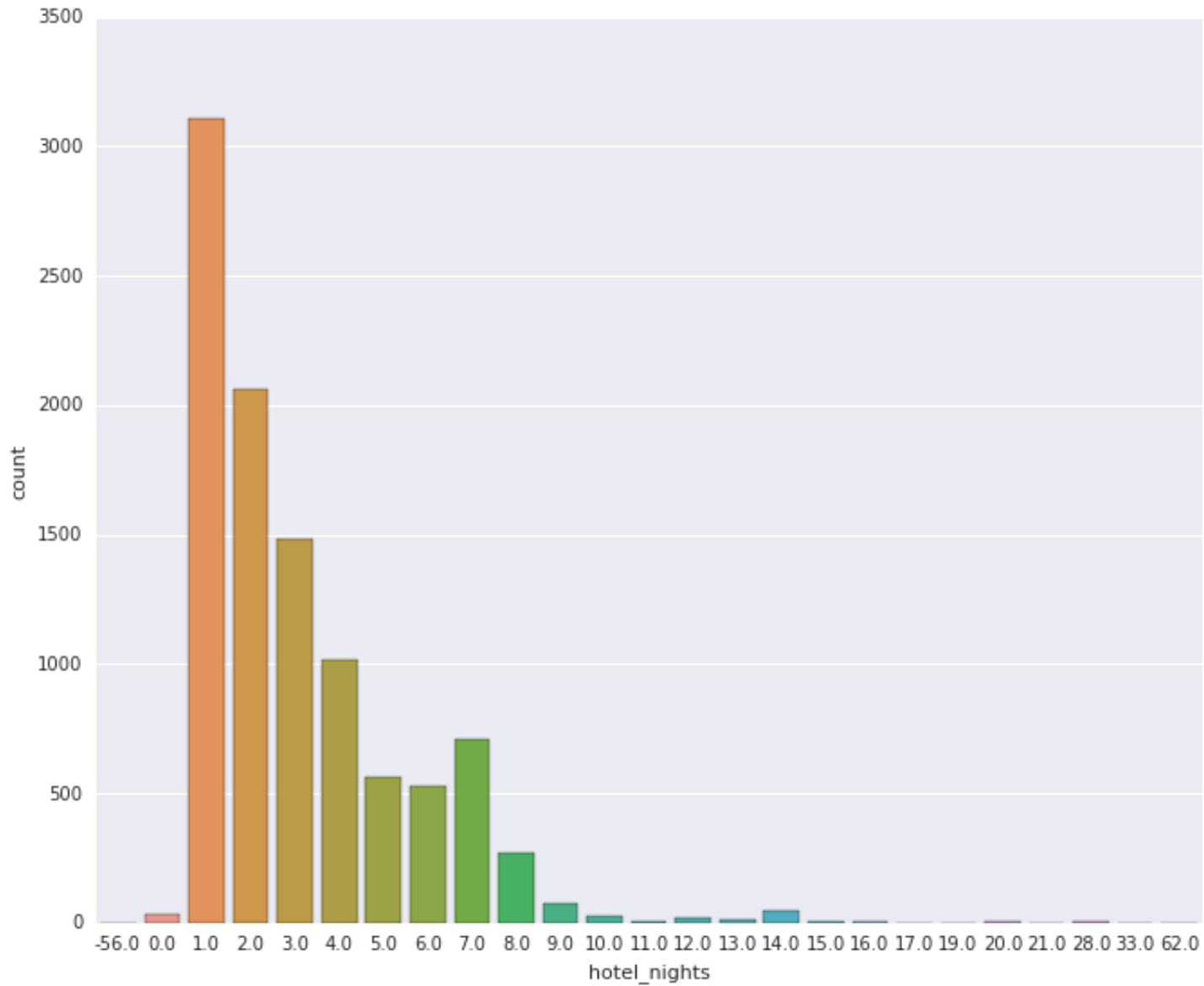
# Expedia – examining features

---

- What are the most countries the customer travel from/to?

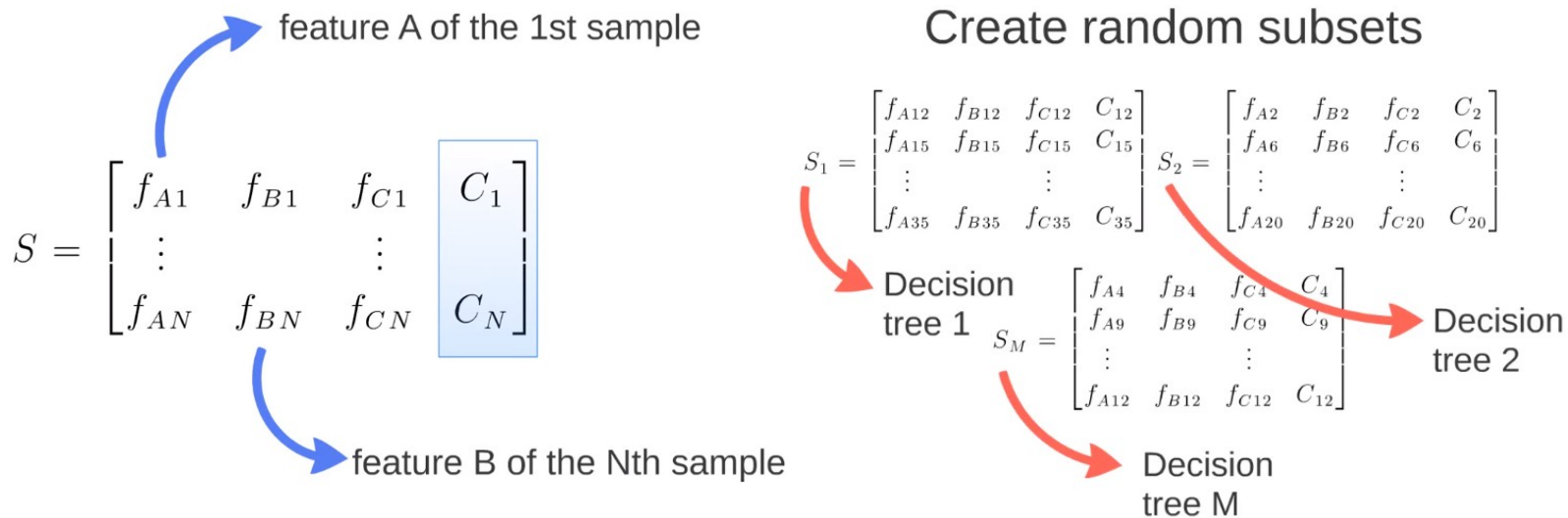


- Nights of stay



# Random Forest Classifier

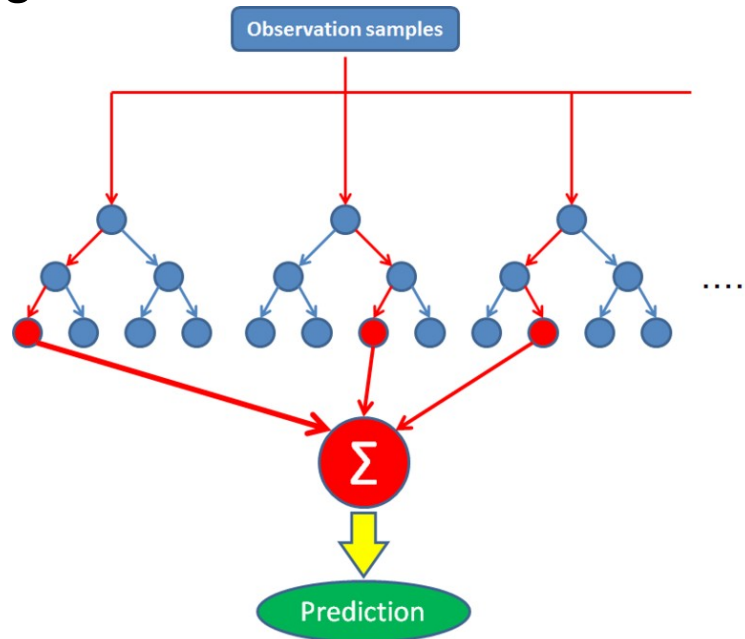
- Supervised learning classifier – Uses bagging methods.
- Random sub-samples.
- Generates decision trees on each sub-sample.



# Random Forest Classifier

---

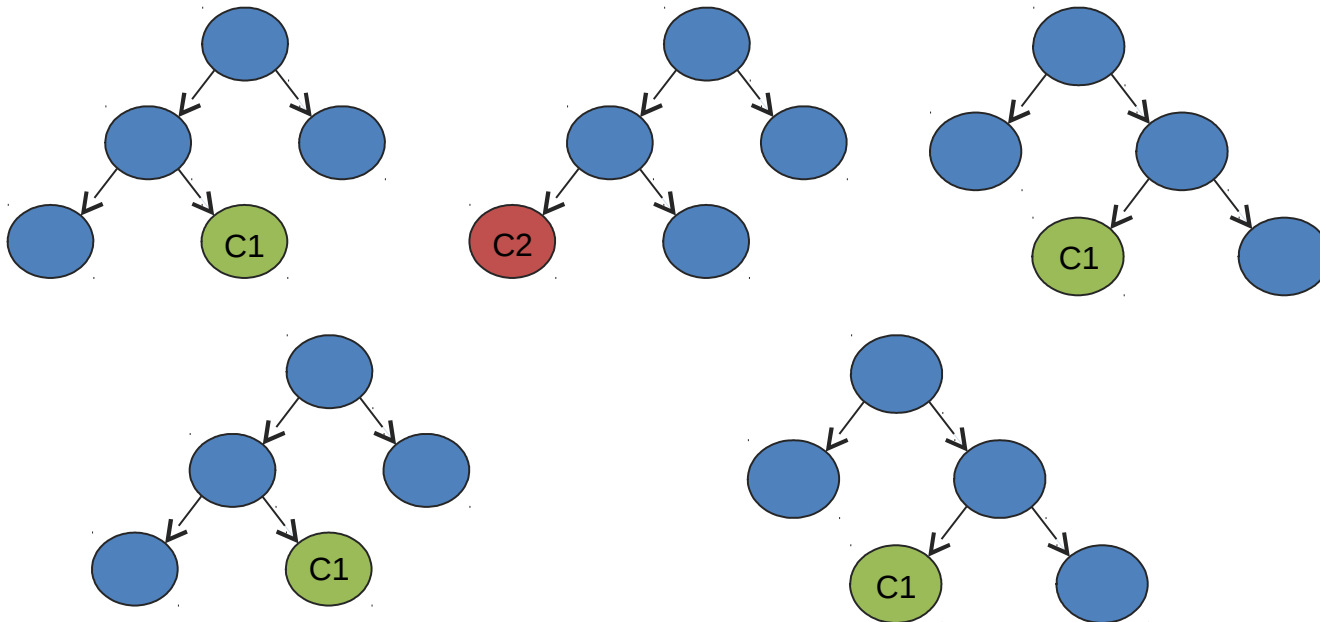
- Sum all the decision trees.
- Mistakes are taken care of.
- The classifier corrects decision trees habit of overfitting to their training set.



# Random Forest Classifier

---

- Why does Random Forest work?
  1. Most trees provide correct prediction for the most part of the data.
  2. Trees make mistake at different place.



# Expedia – How good is the classifier?

---

- We predict 5 hotel clusters for each sample in test.csv
- The evaluation function is **Mean Average Precision @ 5**

$$MAP@5 = \frac{1}{|U|} \sum_{u=1}^{|U|} \sum_{k=1}^{\min(5,n)} P(k)$$

where  $|U|$  is the number of user events,  $P(k)$  is the precision at cutoff  $k$ ,  $n$  is the number of predicted hotel clusters.

Test0 : Truth is 1, Predicted [1,2,3,4,5] => Average precision = 1

Test1 : Truth is 2, Predicted [1,2,3,4,5] => Average precision =  $\frac{1}{2}$       Mean average precision = 0.425

Test2 : Truth is 5, Predicted [1,2,3,4,5] => Average precision =  $\frac{1}{5}$   
Test2 : Truth is 5, Predicted [1,2,3,4,5] => Average precision =  $\frac{1}{5}$

Test3 : Truth is 6, Predicted [1,2,3,4,5] => Average precision = 0  
Test3 : Truth is 6, Predicted [1,2,3,4,5] => Average precision = 0

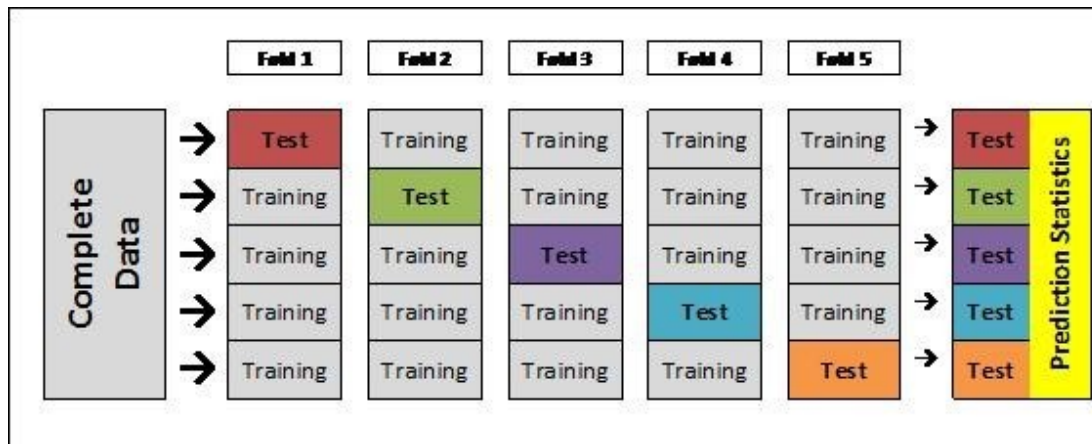




# Expedia – How good is the classifier?

---

- k-fold cross-validation for model tuning



- We could more easily tune the model with a Grid Search for the best parameters



# Expedia - Results

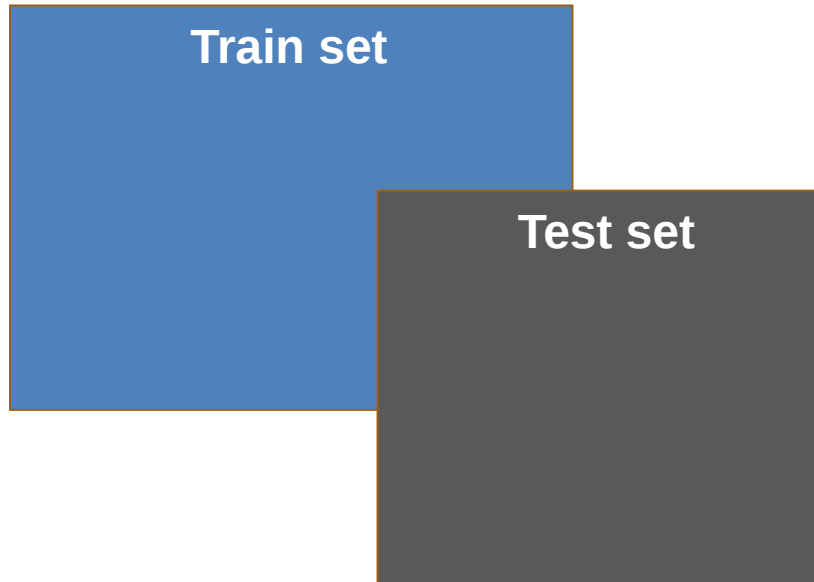
---

- Results with Random Forest classifier:  
0.18584
- Results with most popular local hotels:  
0.30090



# Leakage

---



- user\_location\_country, user\_location\_region, user\_location\_city, hotel\_market and orig\_destination\_distance



# Leakage - Results

---

- Using a more advanced approach with most popular hotels and leakage we got:

0.50050



# Expedia - Conclusion

---

- Machine learning can be used in real-life situations to optimize a product or service
- It is very important to not leak training examples into the test set because the model will overfit
- Here the best model will have to find the leak (1/3) and train itself to catch the rest of the holdout data (2/3)





**LUND**  
**UNIVERSITY**