# Appraiser : How Airbnb Generates Complex Models in Spark for Demand Prediction

hector.yee@airbnb.com

@eigenhector

# Hector Yee
## Selected publications
## 8 Movies & games (Shrek2, Star Wars, C&C etc)



Google Image Search (4 patents)

- Image classification, YH Yee et al, US Patent 8,478,052

- Customizing image search for user attributes, YH Yee, CJ Rosenberg,US Patent 8,782,029

Google Self-driving Car - Perception (1 patent)

Youtube Technical Emmy 2014 (4 papers)

- Label partitioning for sublinear ranking, J Weston, A Makadia, H Yee, ICML 2013

- Affinity weighted embedding, J Weston, R Weiss, H Yee, ICML 2014

# Price tips

# Pricing in a Two-Sided Marketplace

- **Goal:** Equip hosts with a tool to make better informed pricing decisions to meet their needs

- Price tips only

# Overview

- **Modeling**

- **Aerosolve:** open-source ML stack

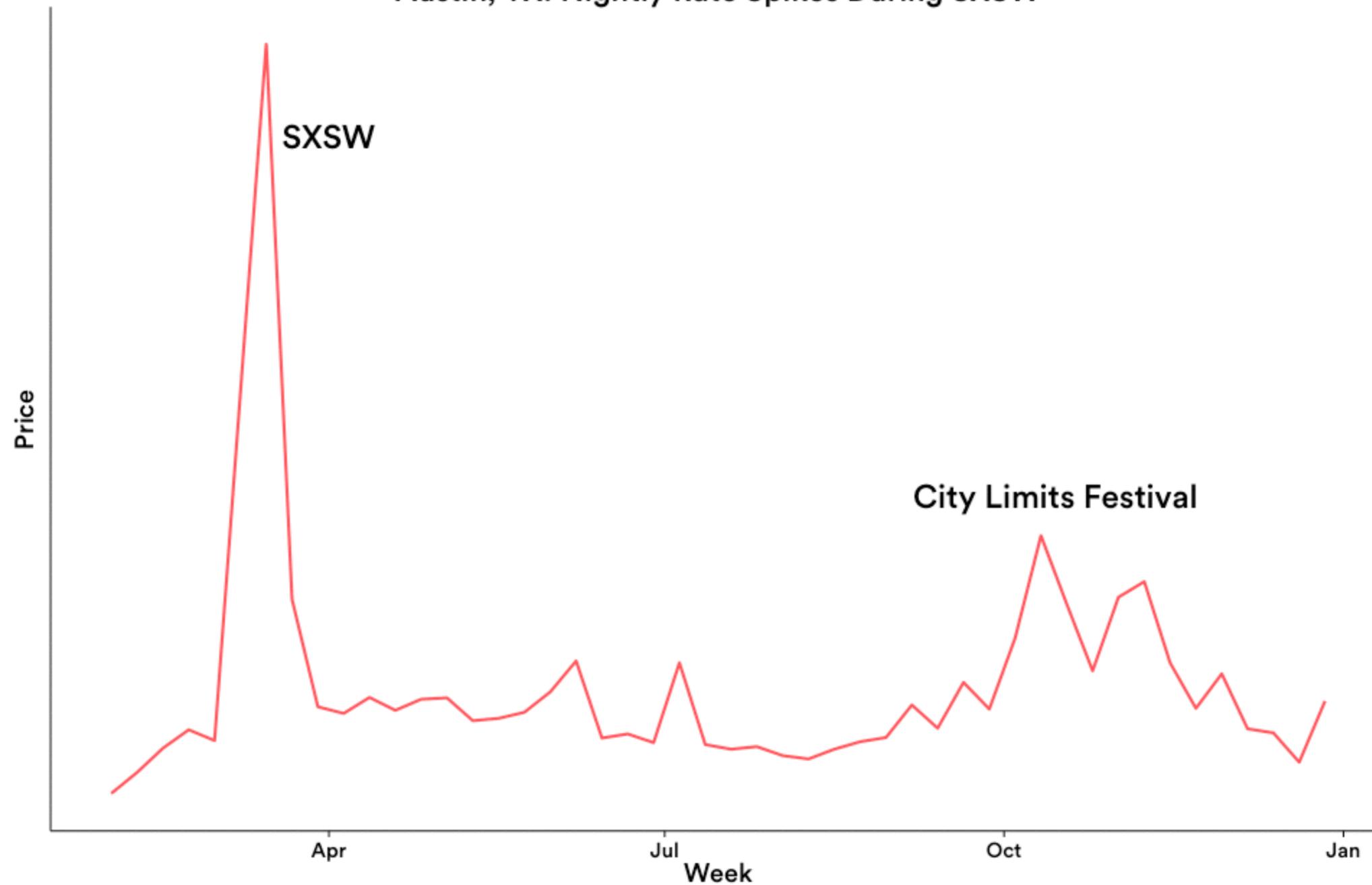# Modeling

# Modeling Approach and Scale

- **Modeling goal: Predict the probability of a booking for each price**

- Scale

  - O(Million) of derived features

  - Over 5B training data points. Proportional to (# listings) x (# days)

# What Affects Prices?

- **Demand** (seasonality, events)
- **Listing location** (market, neighborhood, street block)
- **Listing type and quality**

# Seasonality & Events

Austin, TX: Nightly Rate Spikes During SXSW

SXSW

City Limits Festival

Price

Apr     Jul     Oct     Jan

Week

# Seasonality & Events

**Demand captured by multiple features**



Austin, TX: Model Detects High Demand

# Listing Location

# Listing Location

**Grids / kdtree captures value of locations in San Francisco**
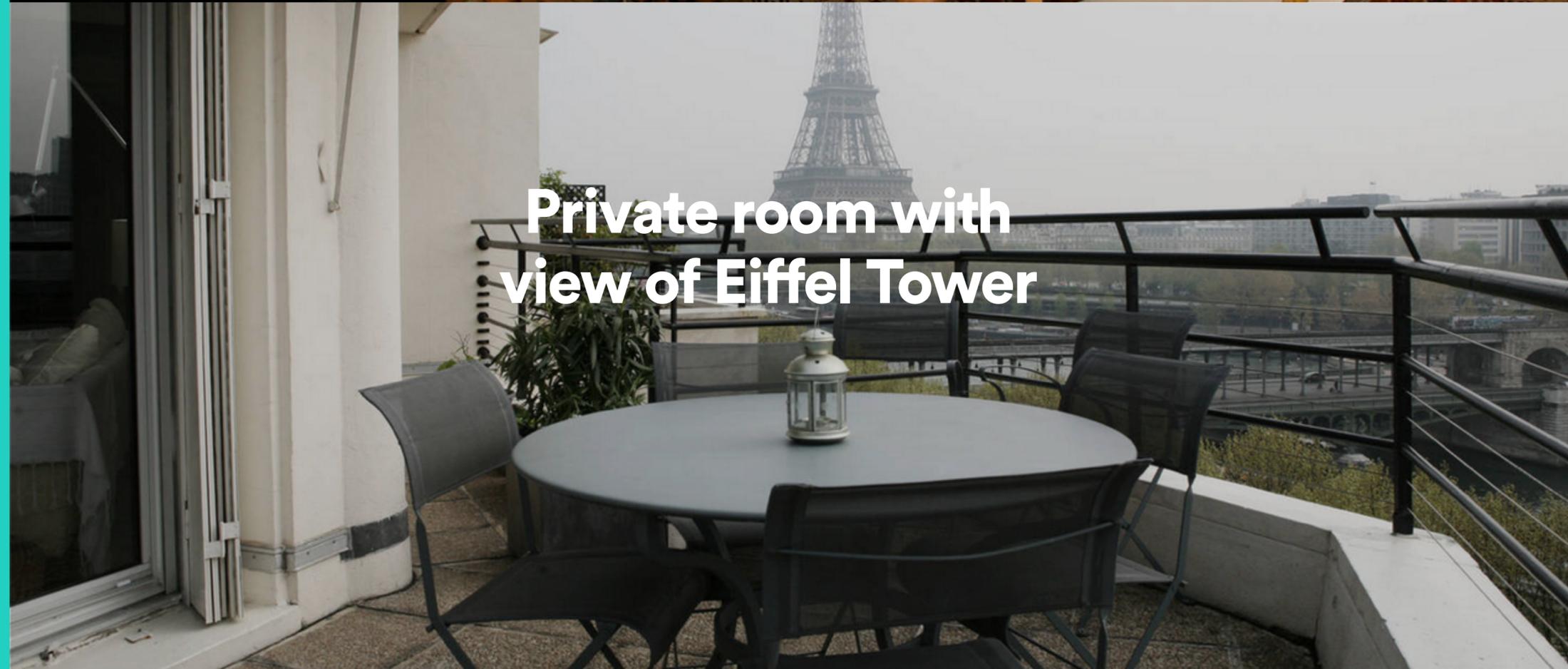
# Listing
# Type & Quality

**Entire apt or a room**

**Amenities**

**Guest reviews**

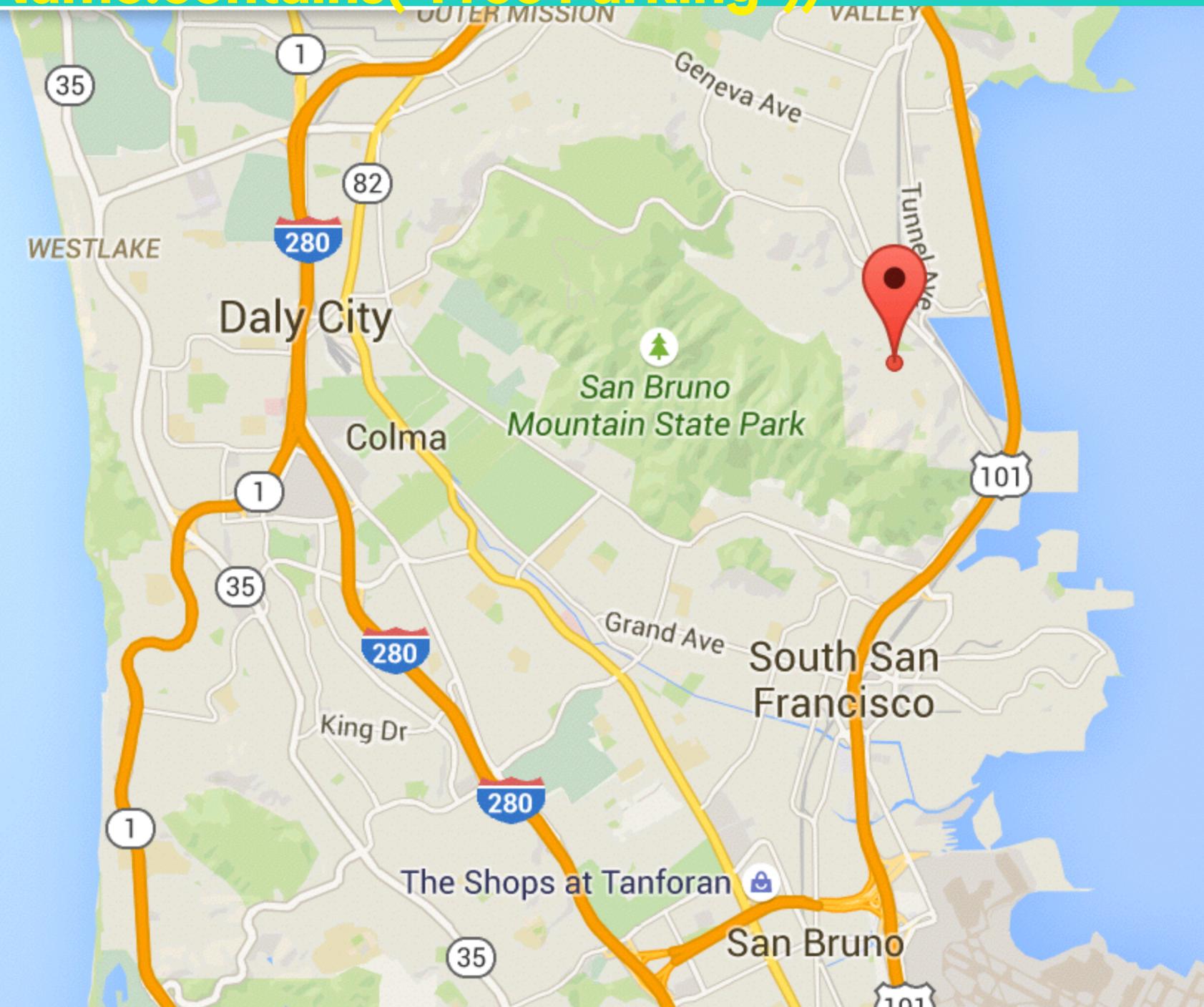Entire home/ houseboat across from Eiffel Tower

Private room with view of Eiffel Tower

# Aerosolve : Machine Learning *for humans*

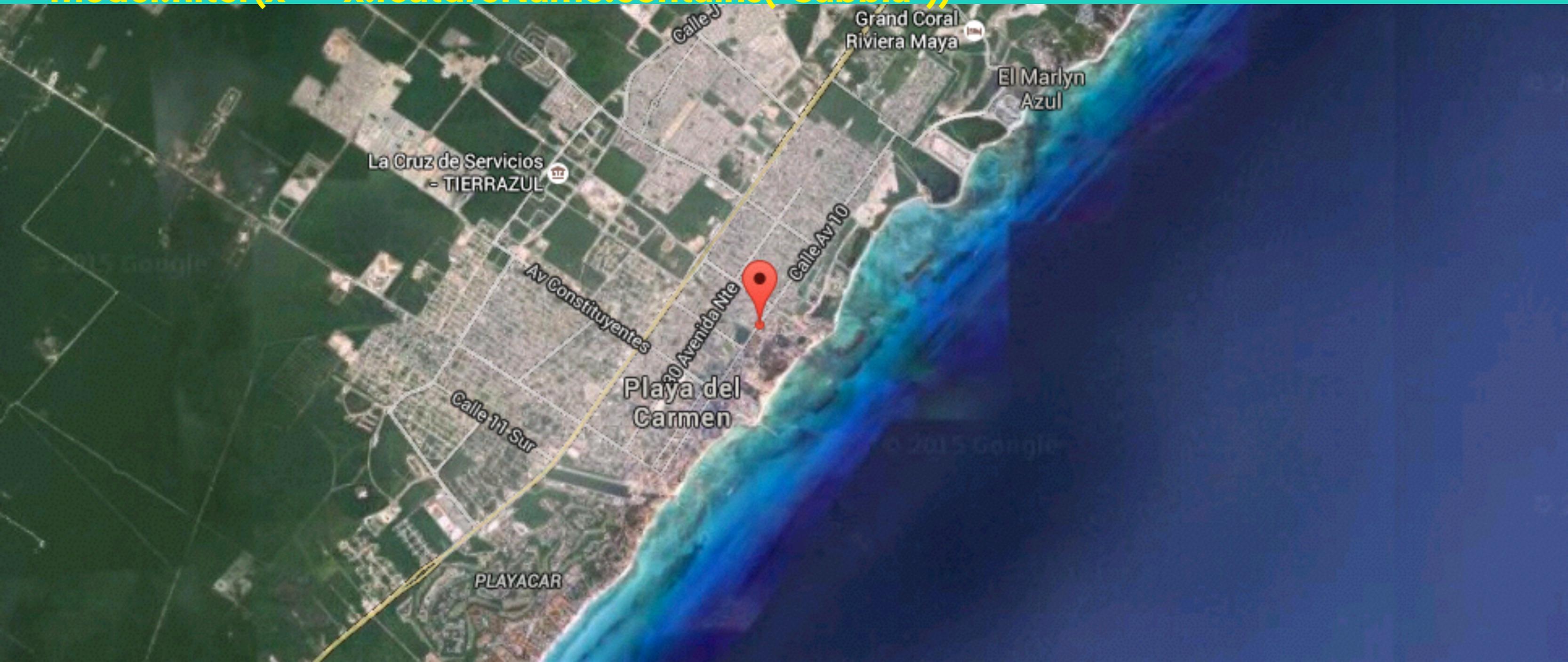# Free parking near SFO Airport?
(37.68359375,-122.40234375

model.filter(x => x.featureName.contains("Free Parking"))

# Sabbia (sand Italian)
## (20.63671875,-87.06640625)

model.filter(x => x.featureName.contains("Sabbia"))

Machine Learning for Humans

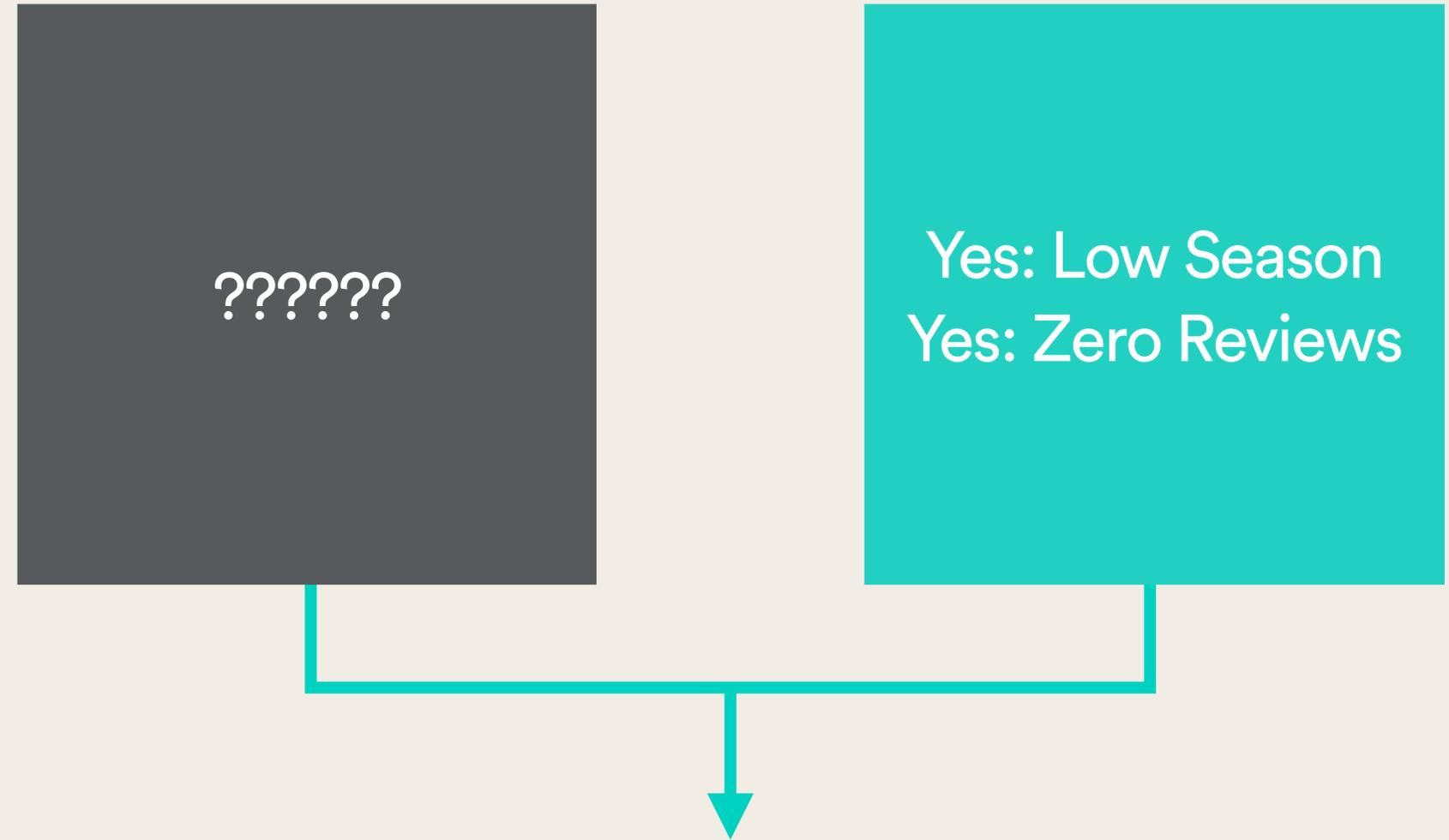Interpretability

"Is my listing priced too high?"

Black Box vs. Glass Box

?????

Yes: Low Season
Yes: Zero Reviews

Answer: Yes
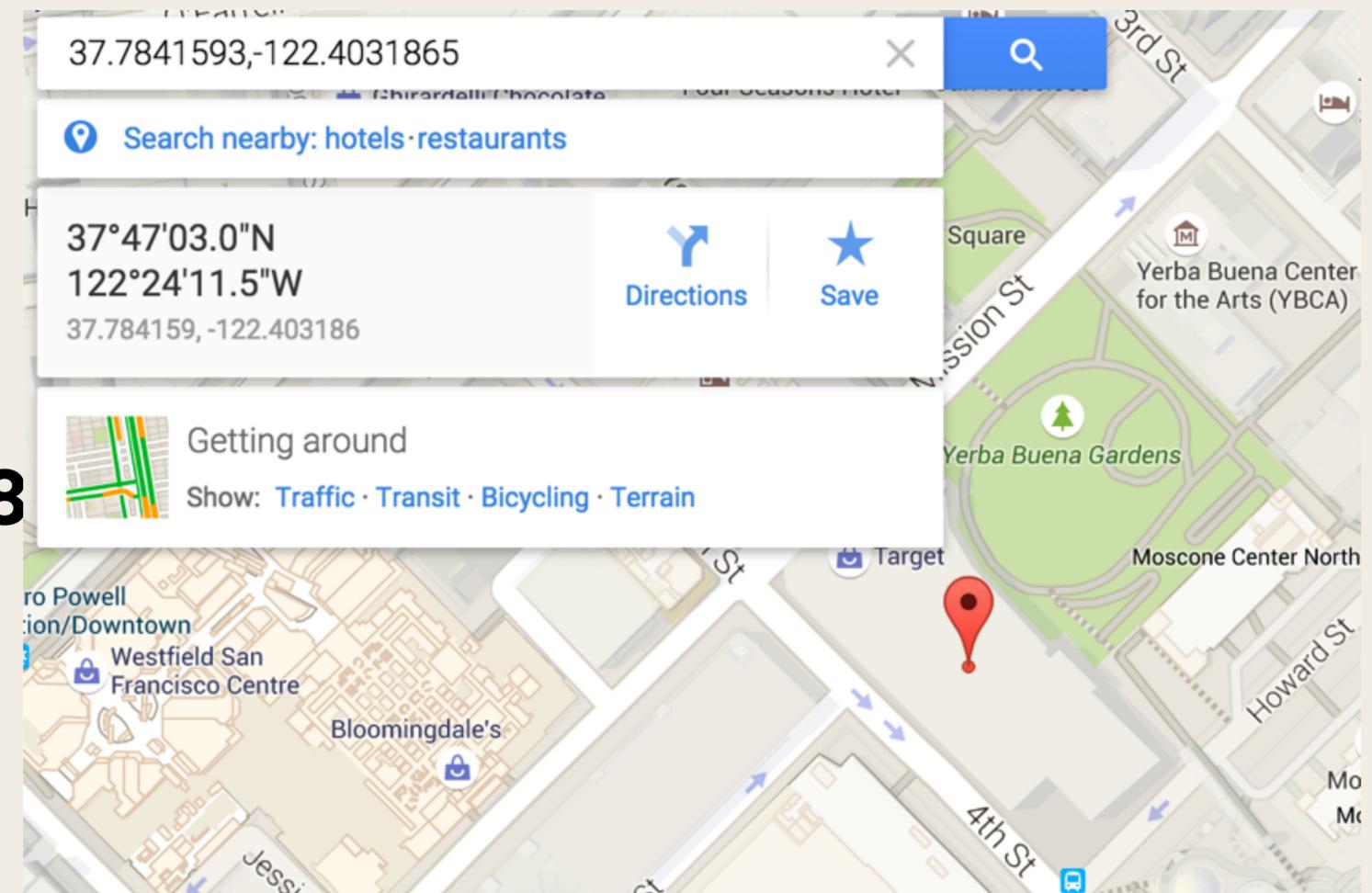
# Question : Will I get a booking at this price?

- Latitude
- Longitude
- Price

**Metreon**

**Lat :37.7841593**

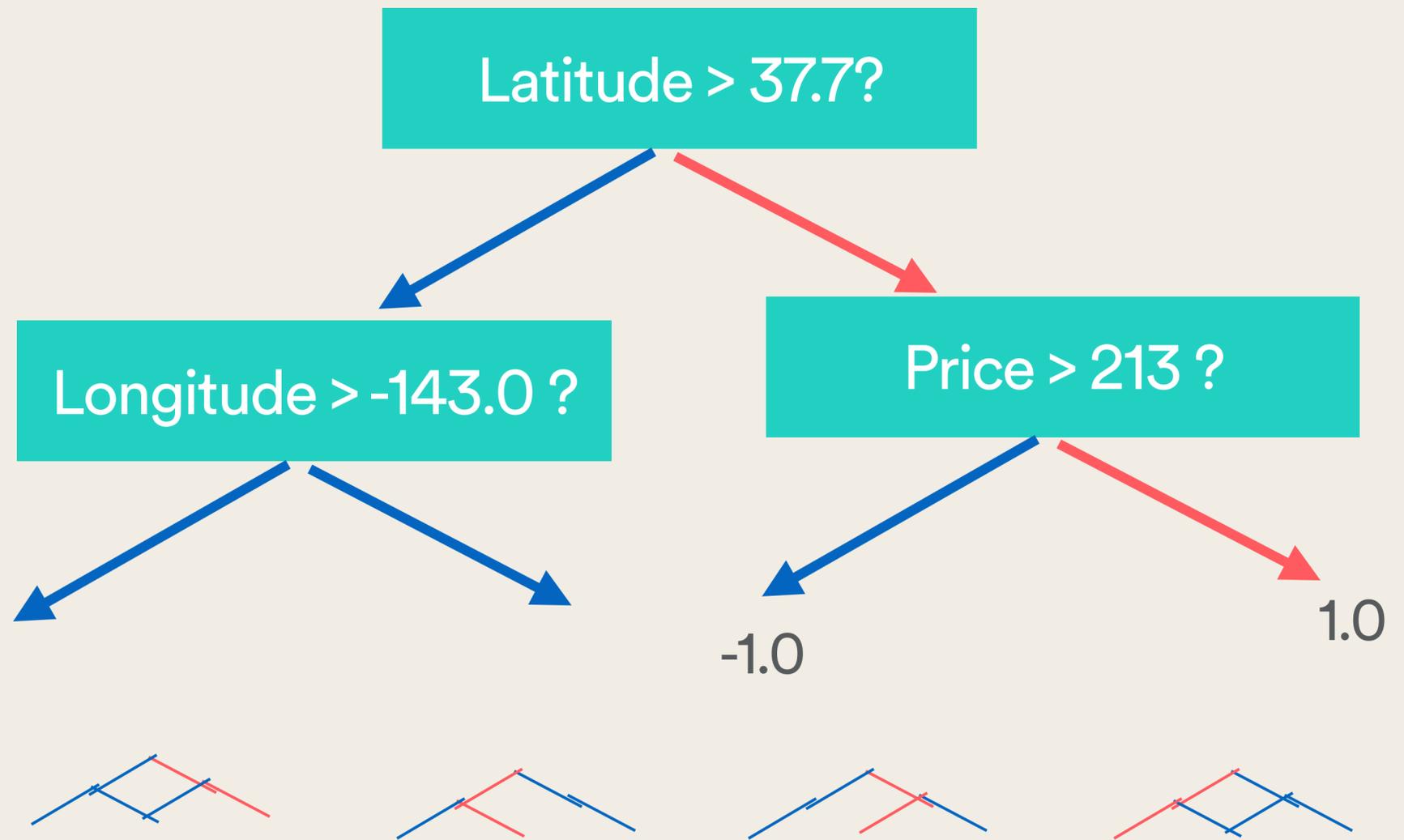**Long : -122.40318**

**Price : 500**

**Black box**

**Random forest**

**Q: Will I get a book?**

**Lat :37.7841593**

**Long : -122.4031865**

**Price : 500**

Latitude > 37.7?

Longitude > -143.0 ?

Price > 213 ?

-1.0

1.0

- Difficult to interpret
  - 1000s of trees in a forest
  - Not clear relationship between variables and target
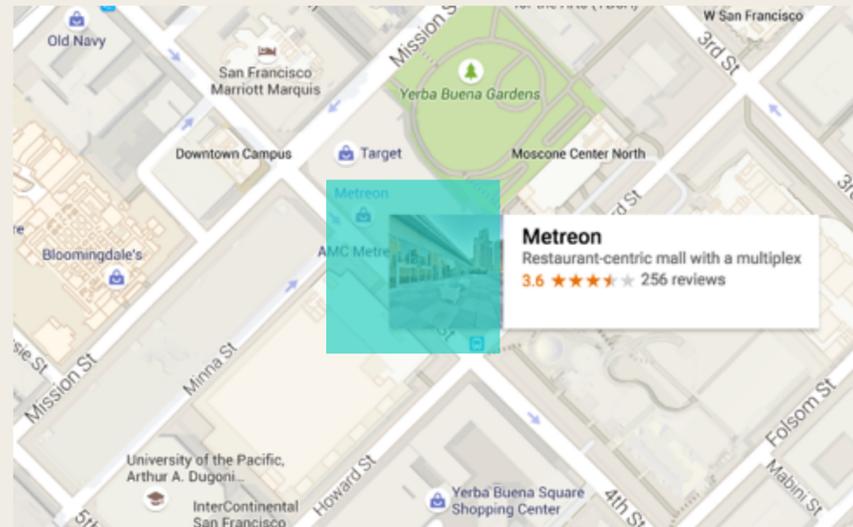
## Glass box model

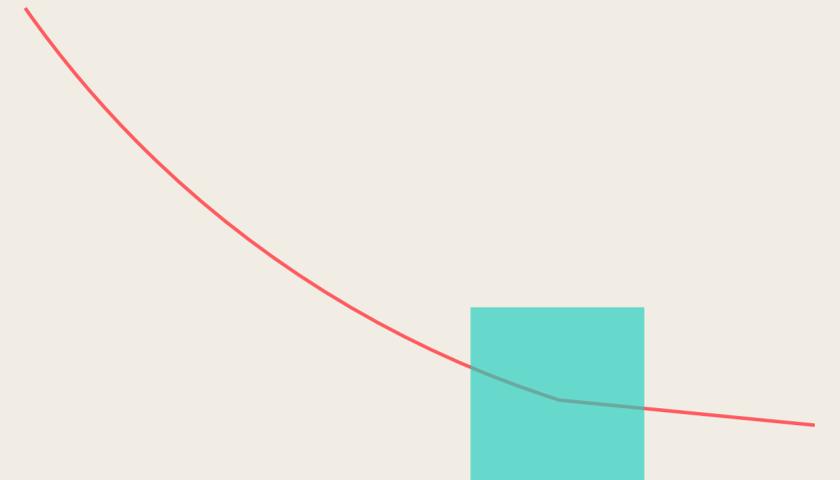**Lat : 37.7841593**

**Long : -122.4031865**

**Price : 500**

- Control quantization
- Control interaction (crosses)
- O(millions) of sparse parameters
- O(tens) active any time

(37.7, -122.4)  AND  Price in [500, 550]



Good location

High price!

# Feature transforms
## Control feature engineering

```
quantize_listing_location {
  transform : multiscale_grid_quantize
  field1: "LOC"
   buckets : [ 0.1, 0.01 ]
  value1 : "Latitude"
  value2 : "Longitude"
  output : "QLOC"
}


quantize_price {
  transform : multiscale_quantize
  field1: "PRICE"
  buckets : [ 10.0, 100.0 ]
  value1 : "$"
   output : "QPRICE" }
```
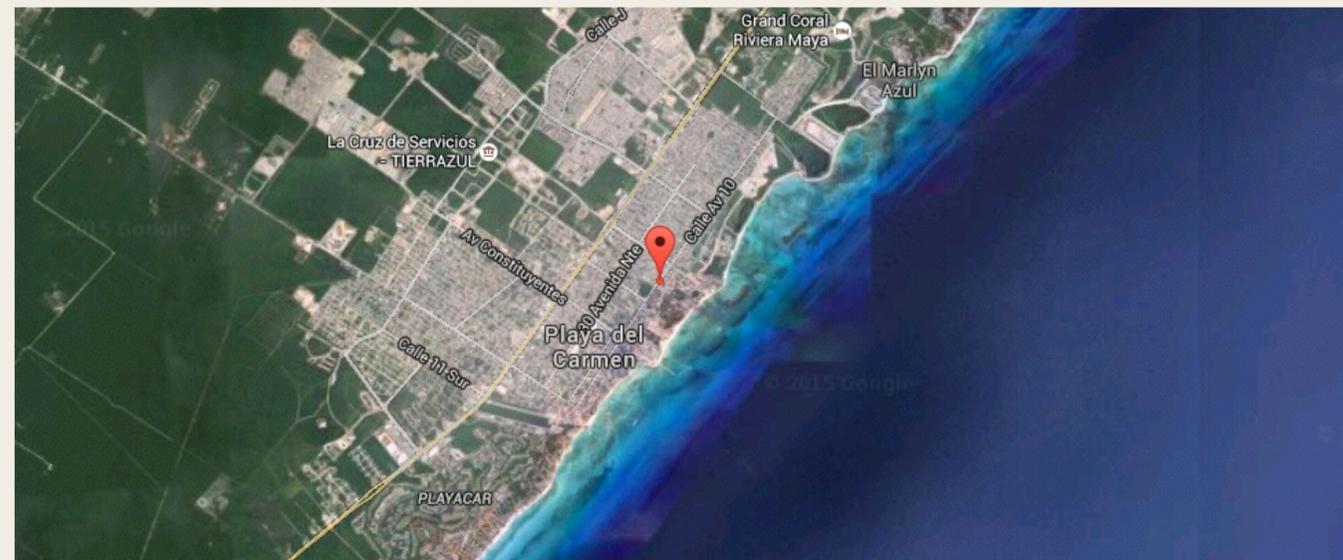
```
Price_X_Location {
  transform : cross
  field1 : "QPRICE"
  field2 : "QLOC"
  output : "PRICE_AND_LOCATION"
}


combined_transform {
  transform : list
  transforms : [
    quantize_listing_location,
    quantize_price,
    Price_X_Location ]
}
```
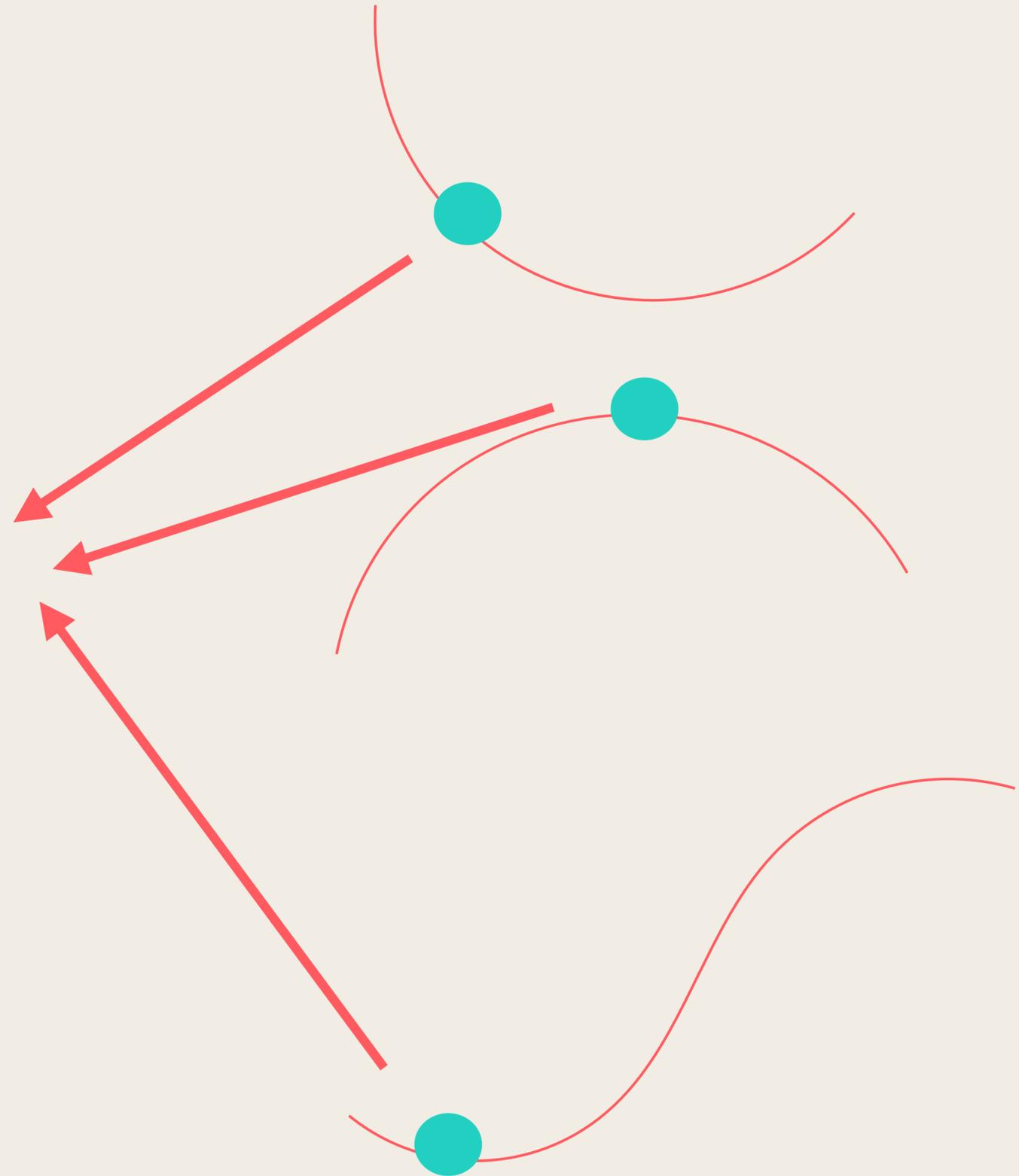
# Benefits of feature transforms

- Write training data once
- Iterate feature transforms on the fly
- Control quantization
- Control interaction (crosses)
- Debuggable models (graphs + readable)
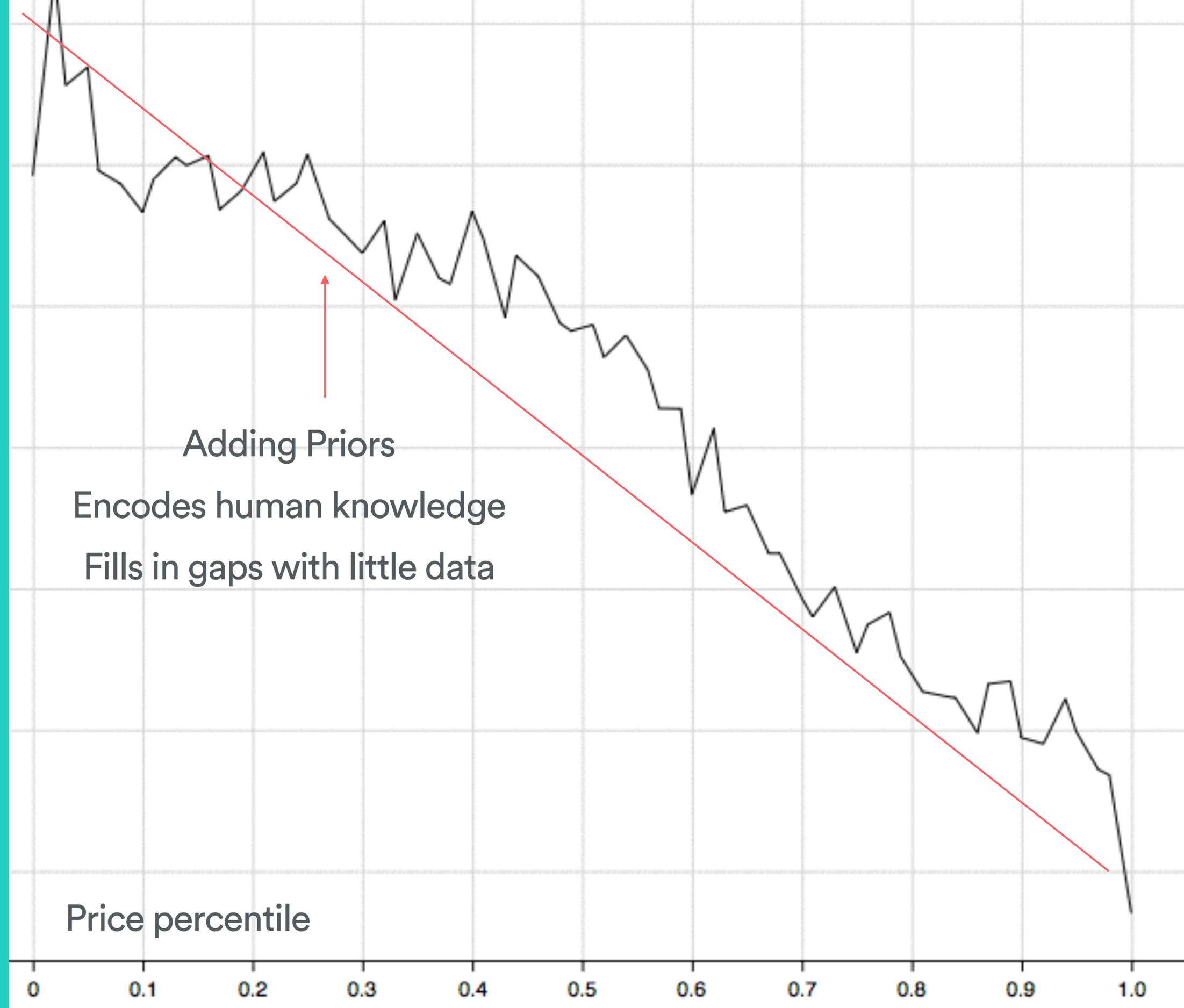
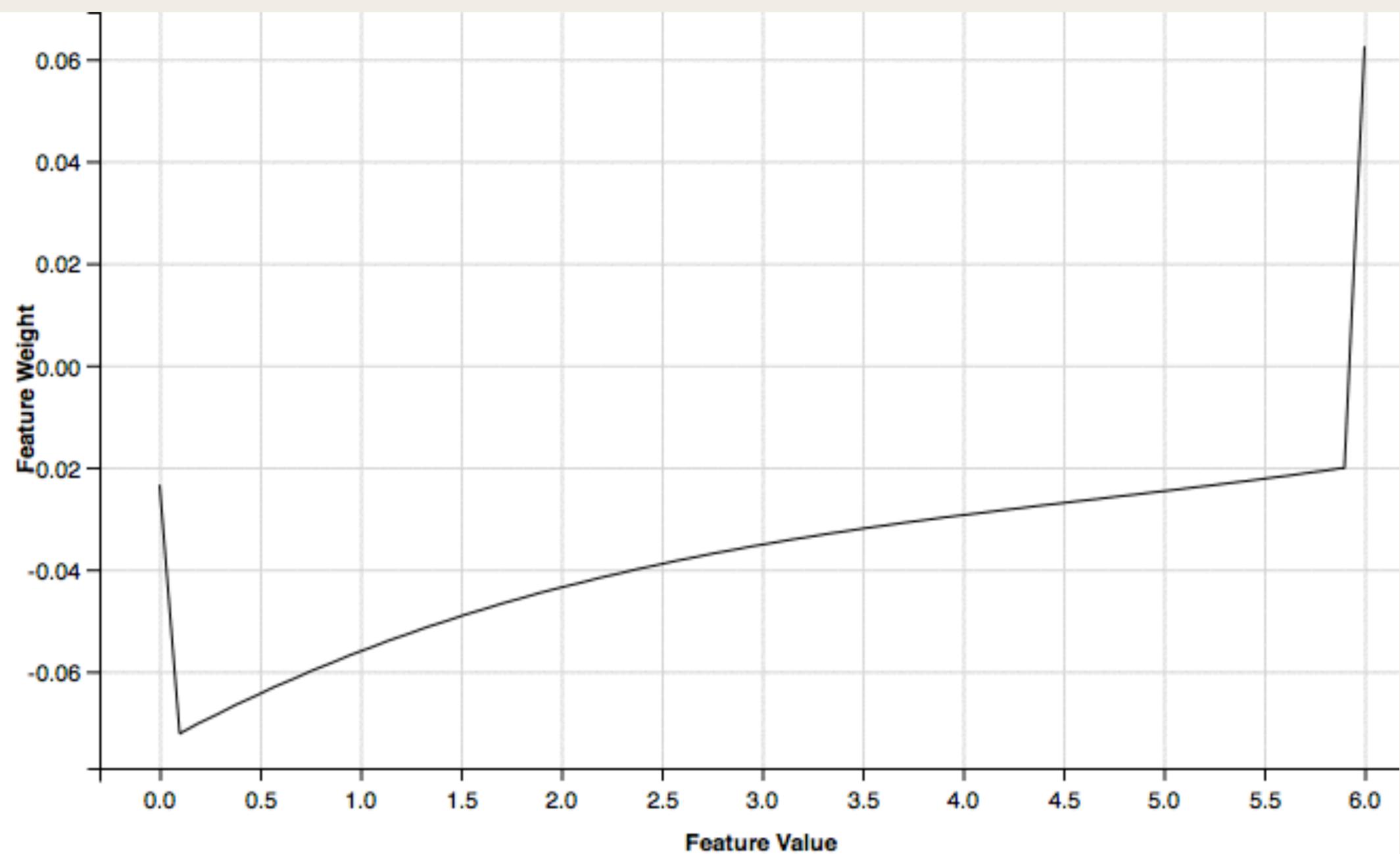**(20.6,-87.0) ^ Sabbia = Good!**

# Additive splines

Weight $= \Sigma$

# Price percentile

**Approximate price percentile in market**

Adding Priors

Encodes human knowledge

Fills in gaps with little data

Price percentile

# Bézier Cubic

## With Dirac Deltas
## End points can vary a lot

# Listing quality

**Dirac + Cubic spline captures effect of reviews**

# reviews

# 3 star reviews

feature_value

0

# reviews

# Debugging Splines

**When splines go bad and weebly wobbly**

High capacity models can cheat!

# Overfitting

**Dangers of high capacity models**

L1 regularization
- Starts learning ID like features
    - Exact location
    - Time of creation
    - Damages Splines

For Linear model:
L1 + L2 + Random Dropout

For Spline model:
L_infinity spline group norm +
Change of basis (Bézier Cubic - 4 params) +
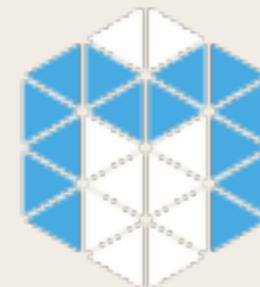Dropout for spline model

# Miscellaneous
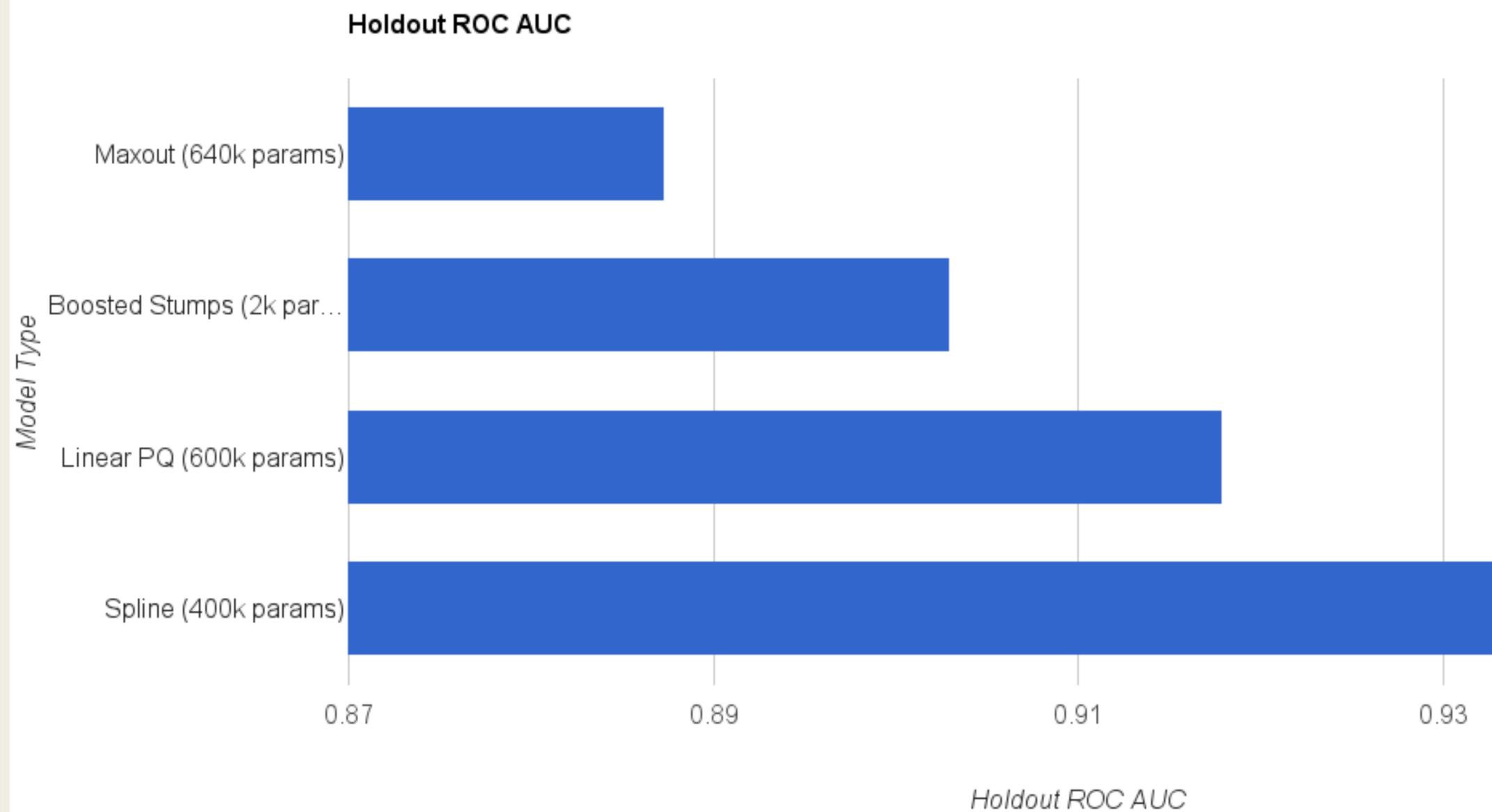
# Achitecture

Airbnb Frontend (Ruby)

Appraiser pricing model (Java)

Aerosolve Machine Learning (Scala)

# KD Trees

**Bodies of water kept separate**

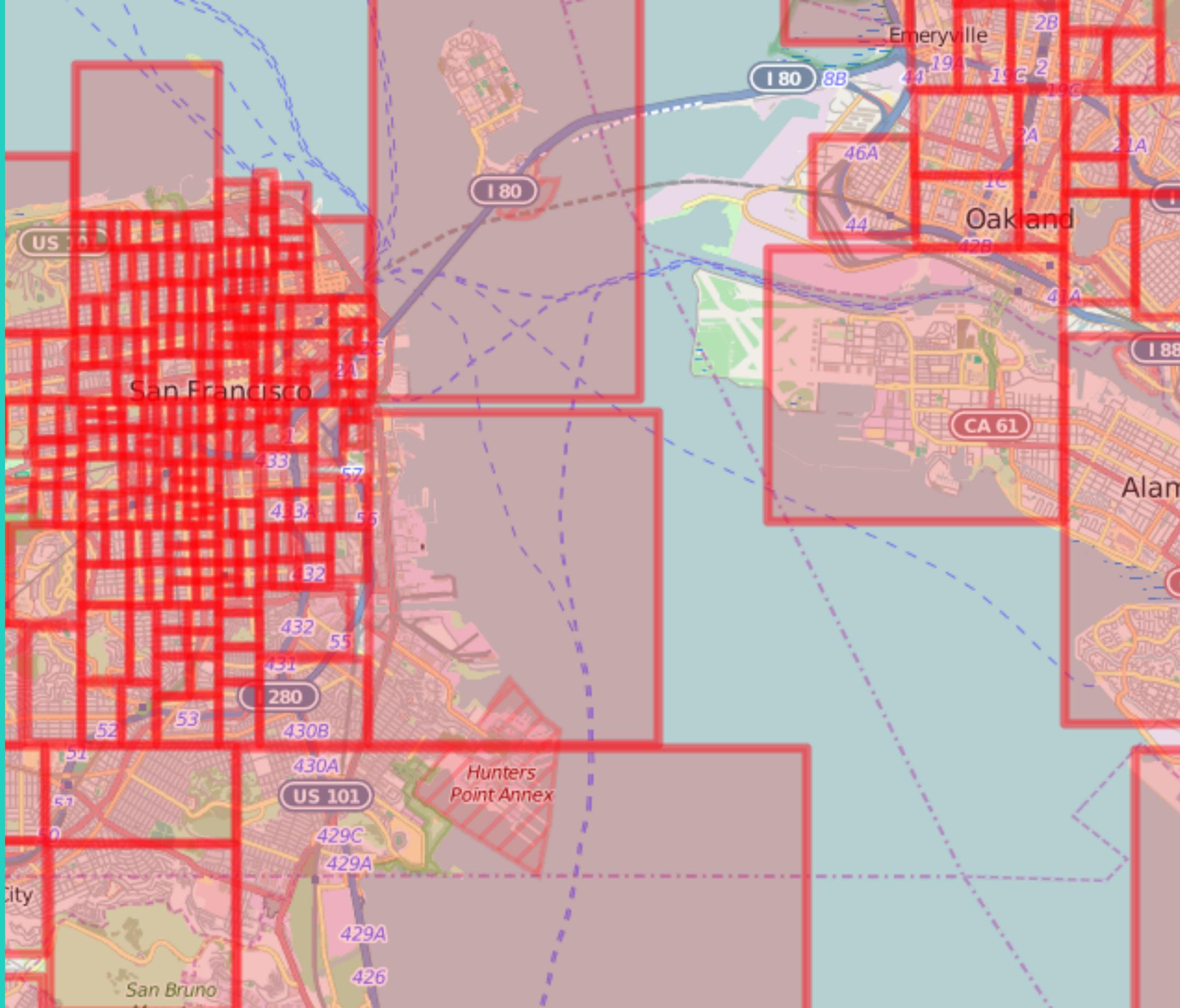**Min Sum $P(leaf_i \mid parent) * n\_leaf_i$**
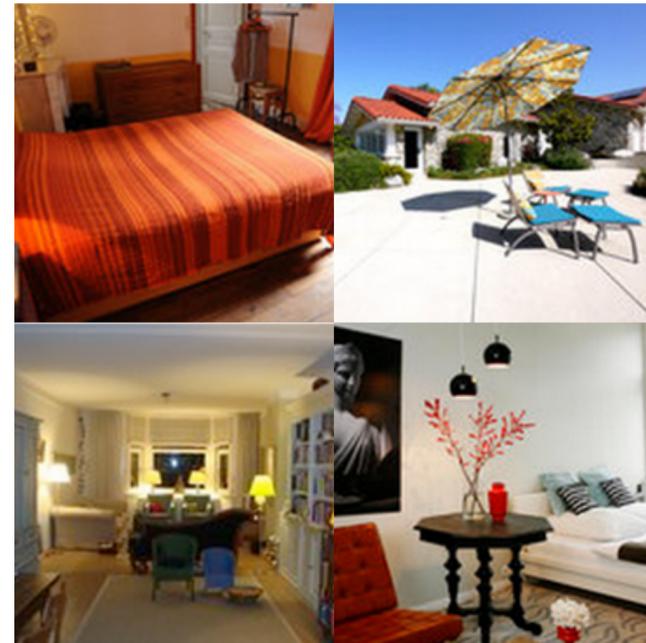
# Image Ranking
## Human curated vs. organic bookings

**Ornate Living Rooms**


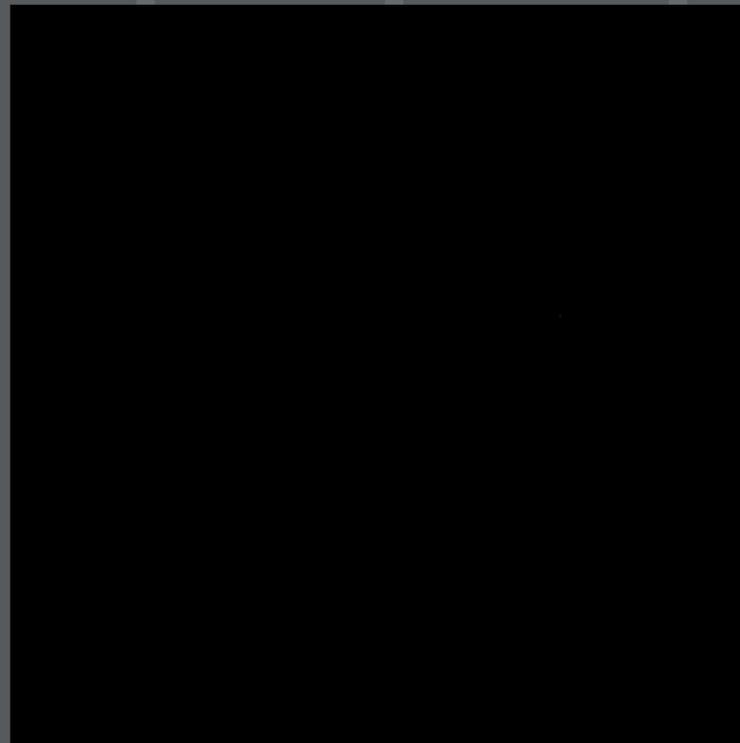
vs.

**Creature Comforts**



### Generate Features

- RGB, HSV, Edge histograms, Texture histograms

### Sparsify

- Winner Take All Hash (max of random elements)

### Ranking Loss

# Demo : Learning to paint

http://airbnb.github.io/aerosolve/

# Questions?

@eigenhector