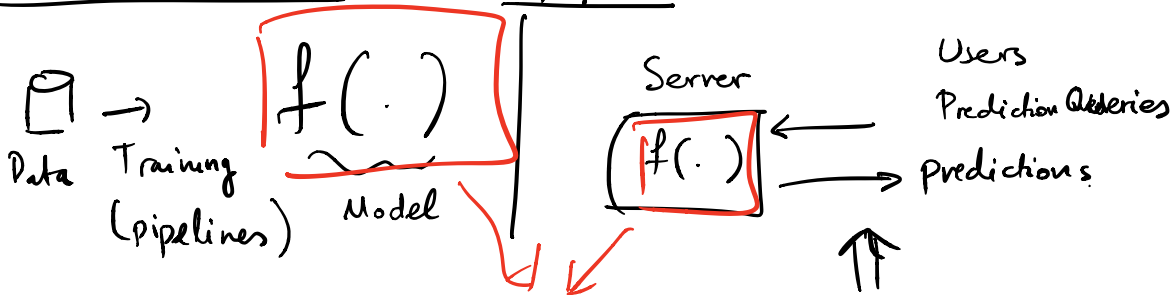


Today: Efficient model serving.

→ Compression, Relational data (embeddings)

Prediction Serving



Serving Systems

Models: as a black box
(reuse the same code)
as training

Standard tricks
for efficient serving.
for user requests.

Problem: Agnostic to ML.

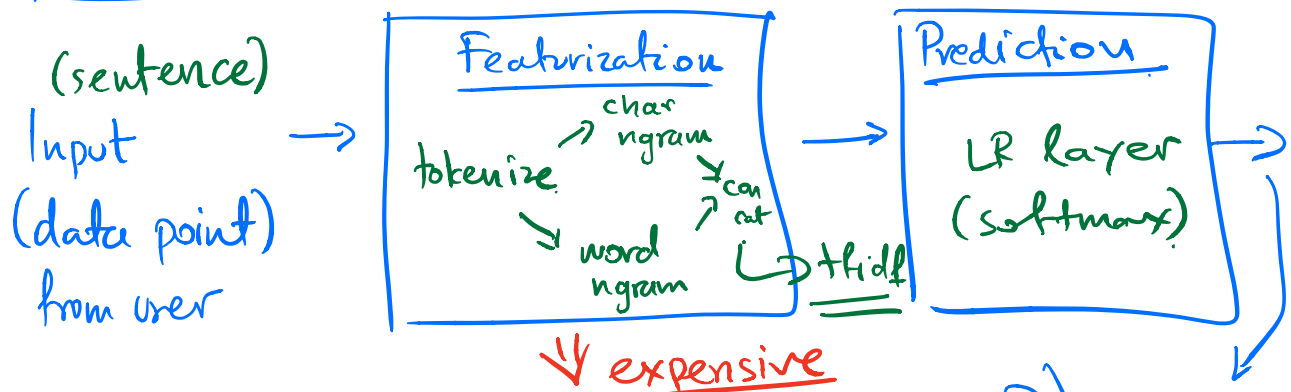
Focus only on prediction serving prob. m.

- 1) User requests
 - * batching
 - * result caching.
- 2) Increase user capacity
 - * replicate the model. (parallel evaluation)

external optimization.

Q: What is the most expensive operation when serving machine learning models?

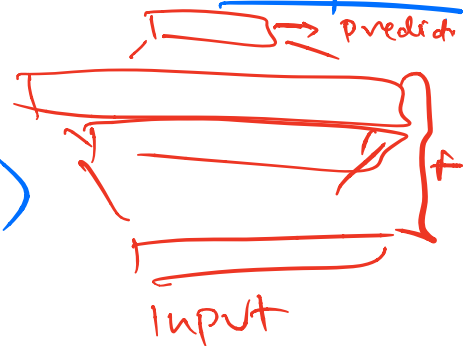
Model computation →



Task: sentiment analysis: (😊 or 😞?) response

Model 1: Logistic Regression

Model 2: LSTM (RNN)



Scenario: Multiple requests from users

Featurization: → split tokens

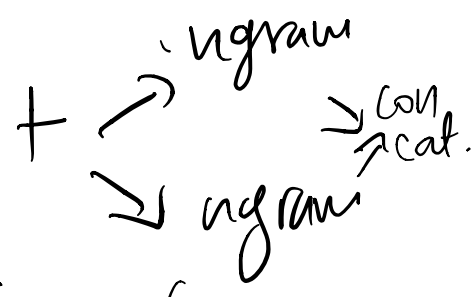
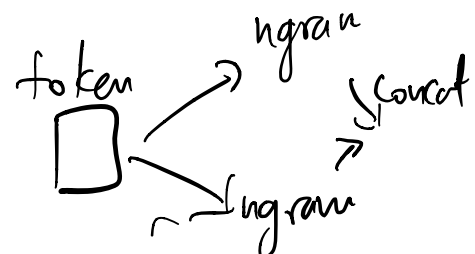
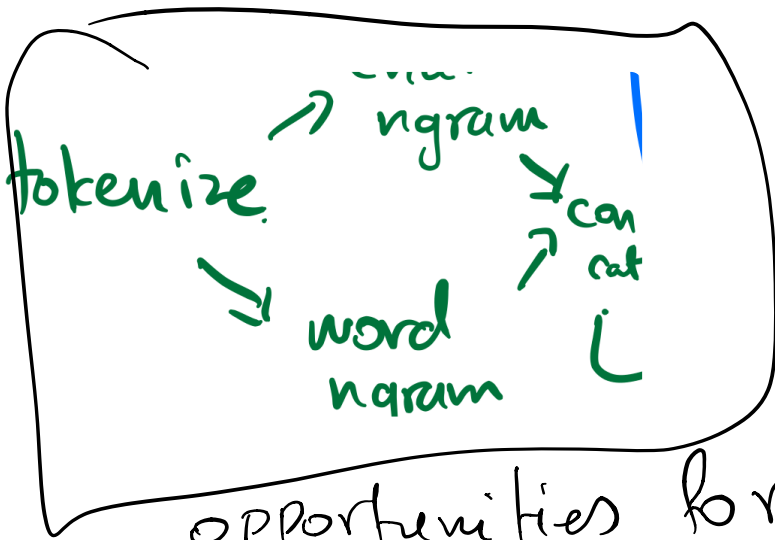
Q: if we consider these ops as a black box?

- extract char n-grams
- word n-grams
- find hidf values
- concat into a vector.

You have two sentences that share 80% of the same words

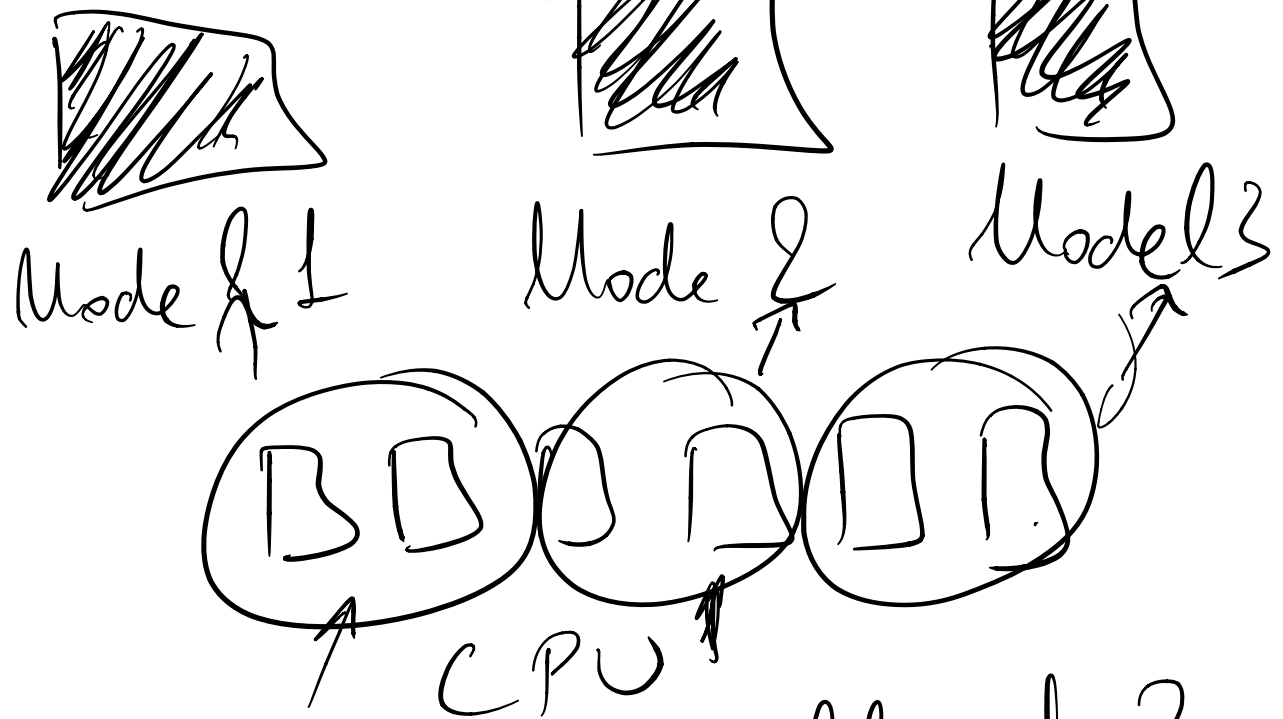
Saving 1: remove redundant lookups. (Hlidf)

Batch → cache my ngrams and compute Hlidf once



opportunities for sharing resources if I open the black box?

Instead of having this
setup



How do you allocate
CPUs to the models?

We need to open the
model black box and look

Thursday

into individual components
and processing?

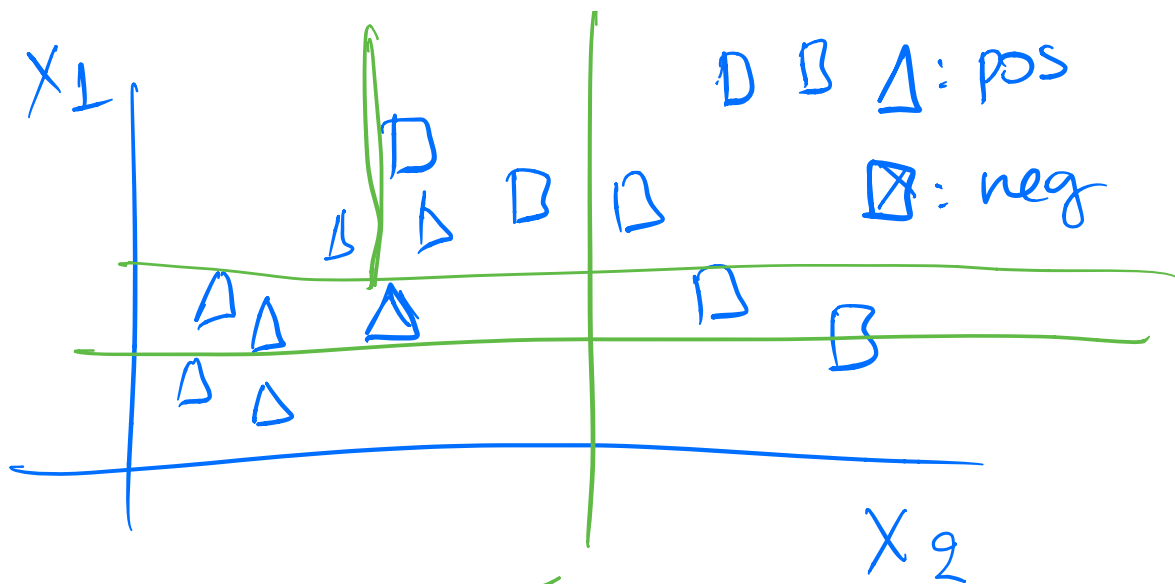
→ No all input are
equally hard.

Model Cascades \leftrightarrow Boosting

↳ weak learners ($p > 0.5$)

Combine weak learners into an
ensemble

the ensemble is a strong
learner ($p > 1 - \epsilon$)



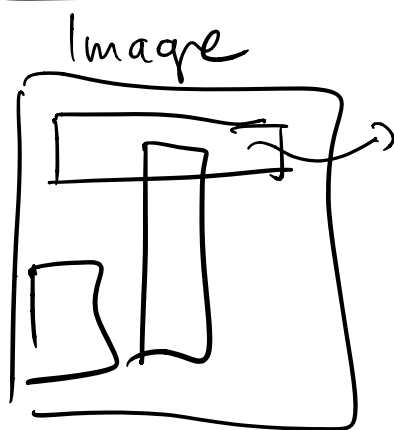
↪ vertical line is a model that uses only 1 out of two features

→ weak learners are much cheaper to evaluate because they rely on fewer features

Cascading classifiers

Viola/Jones Face Detector
2001

Goal: real time face detection



subwindows

start with simple

classifiers that

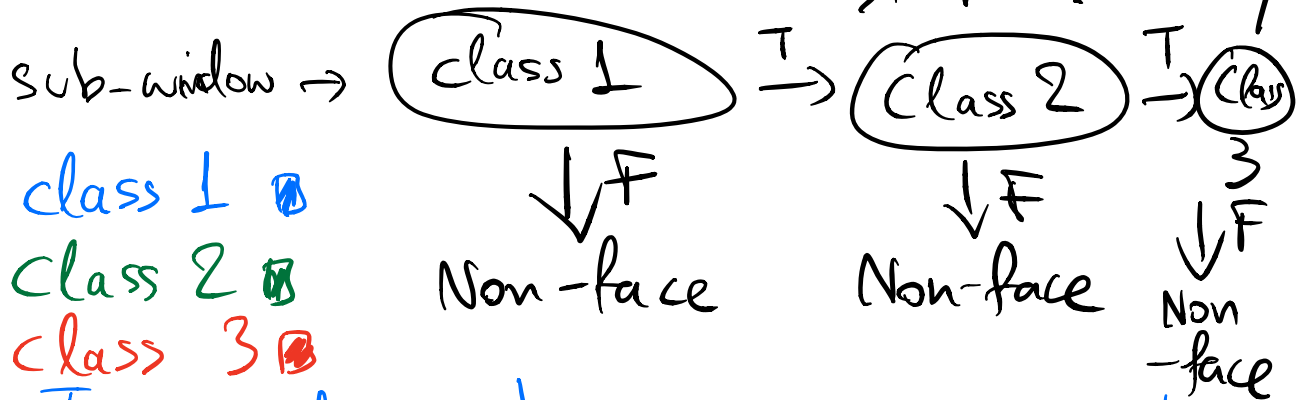
reject many on

the negative

sub-windows while

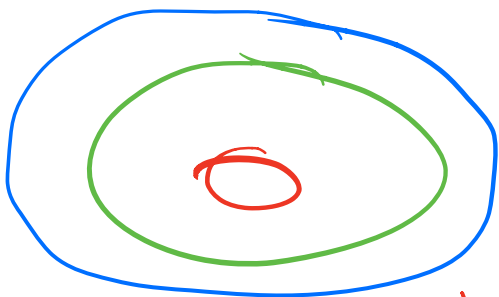
they detect almost all
positive sub-windows

Result from the first classifier → trigger an eval. from the second classifier.



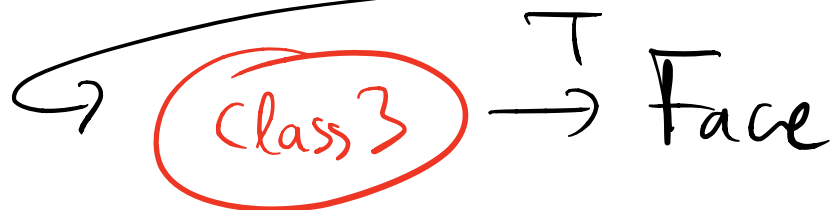
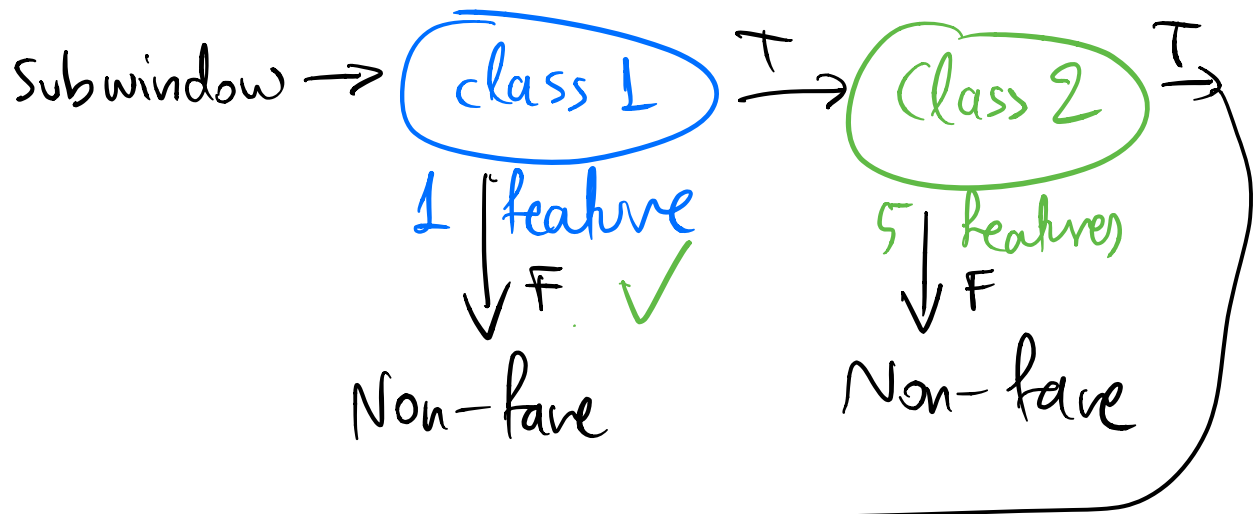
The negative outcome can come out at any point

but the positive outcome has to go through all classifiers



↳ classifiers with lower false positive rate.

↳ class 1 } class 2 } class 3



20 | features
↓ F
Non-face

Goal: filter
negatives
fast

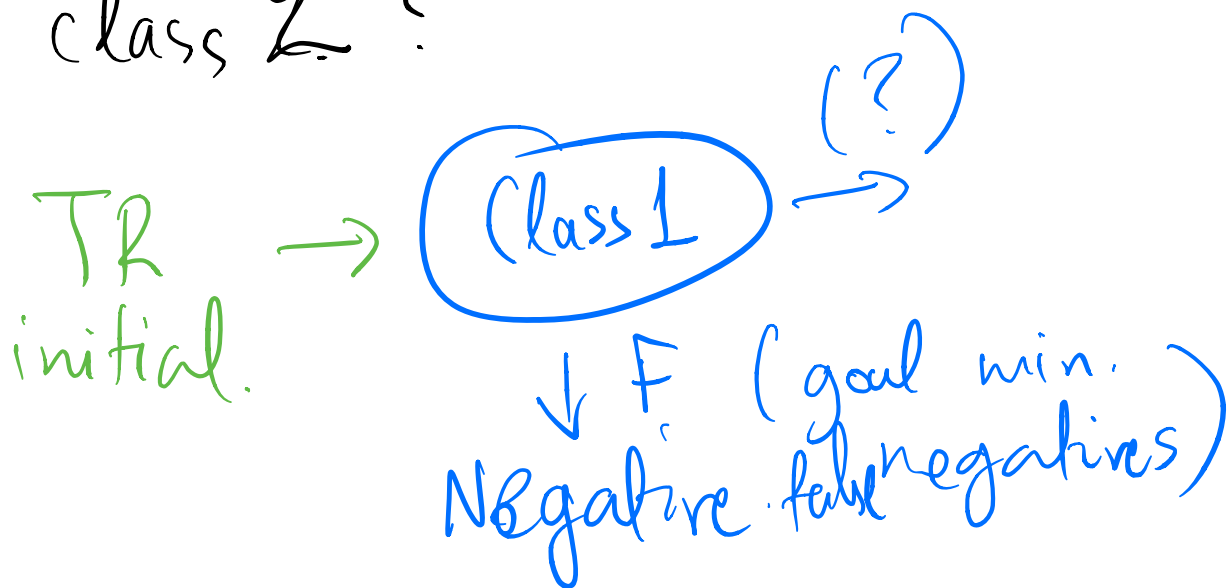
Q: How about training?

A classifier is "optimal" if
it minimizes false negatives

* For each classifier our training goal is to minimize false neg instead of total prediction (classification) error

* Each classifier has a different training set. (TR)

What is the TR of class 2?



TR initial. \rightarrow TR for which
initial Class 1
says "true"

(this dataset still
contains false pos.
non-face elements.)

This training is expensive
because I change the
TR set while I propagate
data through the cascade.

Follow up: questions:

\rightarrow how do I learn the best
cascade so that

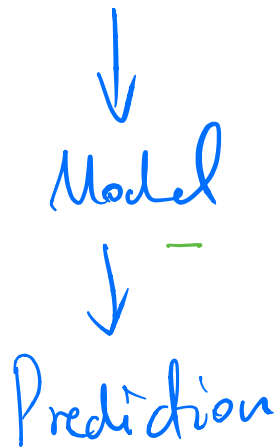
I minimize \rightarrow latency
 \rightarrow resource usage
maximize accuracy.

\rightarrow Willump (MLsys 20)

Statical optimization to correlate features
with different classes

Initial Featurezation Pipeline
(Input to Willump)

Compute all Features



Technical Idea:

Use feature selection

Compute Features

Selected

