

## Linear Interpolation

$$q(w_i | w_{i-2}, w_{i-1}) = \lambda_1 \cdot q_{ML}(w_i | w_{i-2}, w_{i-1}) + \lambda_2 \cdot q_{ML}(w_i | w_{i-1}) + \lambda_3 \cdot q_{ML}(w_i)$$

where  $\lambda_1 + \lambda_2 + \lambda_3 = 1$  and  $\lambda_i \geq 0 \forall i$

(V.V.U.Imp.)

Q: Is  $q(w_i | w_{i-2}, w_{i-1})$  a valid pm estimator?

Ans: Our estimate correctly defines a distribution  
(define  $V' = V \cup \{\text{STOP}\}$ )

$$\sum_{w \in V'} q(w | u, v)$$

$$\begin{aligned} &= \sum_{w \in V'} [\lambda_1 \cdot q_{ML}(w | u, v) + \lambda_2 \cdot q_{ML}(w | v) + \lambda_3 \cdot q_{ML}(w)] \\ &= \lambda_1 \sum_w q_{ML}(w | u, v) + \lambda_2 \sum_w q_{ML}(w | v) + \lambda_3 \cdot \sum_w q_{ML}(w) \\ &= \lambda_1 + \lambda_2 + \lambda_3 \\ &= 1 \quad \text{=: MLE itself correctly defines the distribution} \end{aligned}$$

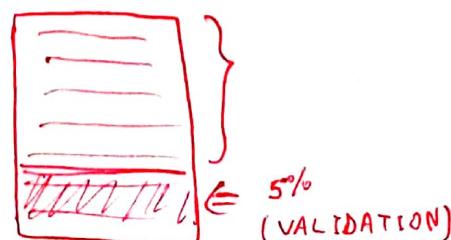
(can also show also that  $q(w | u, v) \geq 0 \forall w \in V'$ )

Trivial

How to estimate the  $\lambda$  values?

- Hold out part of training set as "validation" data.
- Define  $c'(w, w_2, w_3)$  to be the # times the trigram  $(w, w_2, w_3)$  is seen in validation set.
- Choose  $\lambda_1, \lambda_2, \lambda_3$  to maximize:

$$L(\lambda_1, \lambda_2, \lambda_3) = \sum_{w, w_2, w_3} c'(w, w_2, w_3) \log q(w_3 | w_1, w_2)$$



Distributing Indexes : Term Partitioning vs Doc. Partitioning

[V-Tmp. for compre] 5

### Term partitioning

- Dict. of idx terms partitioned into subsets.
  - ↳ Each subset at a node.
  - Along with terms, keep posting lists ↳
- Query routed to terms corresponding to query terms  $\Rightarrow$  greater concurrency

BUT: partitioning by vocab-terms non-trivial (in practice)

- Large cost of merging  $\rightarrow$  Outweighs concurrency
- Another problem: load balancing
  - ↳ governed by distr. of query terms and their co-occurrences
- ★ - makes impl. of dynamic indexing more difficult

### Doc. Partitioning [Better]

- Each node contains idx for a subset of all docs.
- Each query distributed to all nodes — results from various nodes are merged
- More local disk-seeks but less inter-node comm.
- One difficulty: 'idf' must be computed across entire doc. collection even though each node contains a subset of docs. [computed by distributed bg processes]
- How to partition docs to nodes?  
 → Assign all pages from a host to a single node.

13/04/2020

(2)

## Langg. Modelling Appl.

##

- Machine Transl.
- Speech Recog"
- QA / Summarization
- Query Completion

### Imp. Note

#### Langg. Models

Other methods used to improve langg. models:

- o "Topic" or "long-range" features
- o syntactic models

⇒ Linear interpolation is a kind of smoothing.

Another approach :

### Discounting Methods

Define:  $\text{count}^*(x) = \text{count}(x) - 0.5$

"Discounted" counts

Now, estimate =  $\frac{\text{count}^*(x)}{\text{count}(w)}$ , "the"

→ "missing probability mass":  $\alpha(w_{i-1}) = 1 - \sum_w \frac{\text{count}^*(w_{i-1}, w)}{\text{count}(w_{i-1})}$

e.g., in our example,  $\alpha(\text{the}) = 10 \times \frac{0.5}{48} = \frac{5}{48}$

Basic idea:

Distribute these missing mass prob. values to words that were never seen together with 'the' in the training corpus.

### \* Katz Back-Off Models

- $A(w_{i-1}) = \{w : \text{Count}(w_{i-1}, w) > 0\}$
- $B(w_{i-1}) = \{w : \text{Count}(w_{i-1}, w) = 0\}$  → Never seen in training data

- A bigram model

$$q_{BO}(w_i, w_{i-1}) =$$

$$\begin{cases} \frac{\text{count}^*(w_{i-1}, w_i)}{\text{count}(w_{i-1})} & ; \text{if } w_i \in A(w_{i-1}) \\ \alpha(w_{i-1}) \cdot \frac{q_{ML}(w_i)}{\sum_{w \in B(w_{i-1})} q_{ML}(w)} & ; \text{if } w_i \in B(w_{i-1}) \end{cases}$$

In this way  $q_{BO}(w_i, w_{i-1})$  would never be 0.

(3)

### \* Katz Back-Off Model (Trigrams)

$\text{count}(w_{i-2}, w_{i-1}, w_i) - 0.5$

$$\alpha(w_{i-2}, w_{i-1}) = 1 - \sum_{w \in A(w_{i-2}, w_{i-1})} \frac{\text{count}^*(w_{i-2}, w_{i-1}, w)}{\text{count}(w_{i-2}, w_{i-1})}$$

$$A(w_{i-2}, w_{i-1}) = \{w : \text{Count}(w_{i-2}, w_{i-1}, w) > 0\}$$

$$B(w_{i-2}, w_{i-1}) = \{w : \text{Count}(w_{i-2}, w_{i-1}, w) = 0\}$$

\*  $q_{ML}$  of unigram is well-defined since we only consider the words that are not in vocab. (But not for bigram)

∴ we use  $q_{BO}(b_i)$  for  $q_{BO}(\text{tri})$  → Remember

### • A trigram model

$$q_{BO}(w_i | w_{i-2}, w_{i-1}) = \begin{cases} \frac{\text{count}^*(w_{i-2}, w_{i-1}, w_i)}{\text{count}(w_{i-2}, w_{i-1})} & ; \text{ if } w_i \in A(w_{i-2}, w_{i-1}) \\ \alpha(w_{i-2}, w_{i-1}) \cdot \frac{q_{BO}(w_i | w_{i-1})}{\sum_{w \in B(w_{i-2}, w_{i-1})} q_{BO}(w | w_{i-1})} & ; \text{ if } w_i \in B(w_{i-2}, w_{i-1}) \end{cases}$$

### \* LM: Summary

##

3 steps:

- Expand  $p(w_1, w_2, \dots, w_n)$  using chain rule
- Make Markov Independence Assumptions
- smooth the estimates using lower order counts.

## Evaluation of Lang. Models

2

### Perplexity:

- test data:  $\frac{m}{\downarrow}$  sentences :  $s_1, s_2 \dots, s_m$

For each sentence, assign a prob. value

$$- \text{Prob. under model} (\log \text{prob.}) = \log \prod_{i=1}^m p(s_i) = \sum_{i=1}^m \log p(s_i)$$

[This prob. value should be higher]

$$- \text{Perplexity} = 2^{-l}$$

$$\text{where } l = - \frac{1}{M} \sum_{i=1}^m \log p(s_i) \rightarrow \text{Higher if better}$$

Depends on test data size ( $m$ )  
v.v. Imp

Normalizing by  $\frac{M}{\downarrow}$

Total # words in the test data

perplexity  $\rightarrow$  lower if better

### Intuition

$$- \text{Vocab } V: N = |V| + 1 \quad (\text{STOP word})$$

$$\text{Model: } q(w|u,v) = \frac{1}{N} \quad (\text{Uniform prob.})$$

$$\forall w \in V \cup \{\text{STOP}\} \quad \forall u, v \in V \cup \{\#\}$$

$$\text{In this case, perplexity} = 2^{-l} \quad \text{where } l = \log \frac{1}{N}$$

$$= \boxed{N}$$

\* Perplexity is a measure of effective "branching factor"

Lower perplexity value  $\Rightarrow$  More context (5)

↳ Lower in case of trigrams (compared to bigrams and unigrams)

## \* Langg. Modelling Variants

$$\bullet P(x_1, x_2, \dots, x_n) \equiv P(x_m | x_1, x_2, x_3, \dots, x_{n-1}) \Leftrightarrow (\text{chain rule})$$

$$\bullet P(x_i | x_1, x_2, \dots, x_{i-1}, x_{i+1}, x_{i+2}, \dots)$$

Predicting middle value based on left and right

$$\bullet \text{LM at char. level in place of } \underline{\text{words}}.$$

Consider word as a seq. of char.

$\Rightarrow$  Whole sentence is a seq. of char.

## \* State-of-the-art in LM

- Tri, Bi  $\rightarrow$  20 yrs. ago. (Rarely used)

(SOTA)

- Models the whole seq. of words (rather than tri-gram or n-gram)  
(No Markov assumption)

Any word in sentence is dependent on whole seq. of words appearing before it.

- Trained on large datasets ( $\sim 1B$  tokens)
- Around half-a-million parameters

# \* BERT: Bidirectional Encoder Representations from Transformers (6)

- MASK: Word that is missing
- Transformer Encoder: A layer of NN archi.
- Classification layer: Fully-connected layer + GELU + Norm  
(lots of parameters and expensive to train)

## Back to Machine Translation

$p(t)$  - The langg. Model

[Done - Trigrams]

$p(s|t)$  - The Translation Model [Now]

## IBM Model 1: Alignments

### Modelling $p(s|t)$

- problem: Data sparsity

$s_1, s_2, \dots, s_m$  m

\$

$t_1, t_2, \dots, t_k$  k

- o Alignment ' $a$ ' identifies which target word each src word originated from

- There are total  $\boxed{(k+1)^m}$  possible alignments

- Formally, an alignment  $\bar{a}$  is  $\{a_1, \dots, a_m\}$  where each  $a_i \in \{0, \dots, k\}$ .

In IBM model 1, all alignments ' $a$ ' are equally likely:

$$p(a | e, m) = \frac{1}{(k+1)^m}$$

← Major simplifying assumption

15/04/2020

(7)

$p(s|t)$ : translation model      src: French (f)

tgt: English (e)

$$f = f_1 f_2 f_3 \dots f_m \quad e = e_1 e_2 e_3 \dots e_k \Rightarrow p(f|e)$$

Two choices involved:

- ① The length of the French sentence ( $m$ )
- ② The choice of words:  $f_1, \dots, f_m$

Let's assume

Alignment:

- ① Each Fr. word is aligned to exactly 1 En. word
- ② Alignment is many-to-one.
- ③ A French word can be generated by no En. word  
 ↳ Special case:  $e_0$  (NULL)

Now we attempt to model  $p(f_1 \dots f_m | e_1 \dots e_k, \text{NULL})$  using alignment variables.

Joint distribution (Marginal)

$$p(f_1 \dots f_m | e_1 \dots e_k, m) = \sum_{a_1=0}^k \sum_{a_2=0}^k \dots \sum_{a_m=0}^k p(f_1 \dots f_m, a_1 \dots a_m) \underbrace{e_1 \dots e_k, m}_{\text{How to estimate this?}}$$

How to estimate this?

## ⑧ IBM Model 2 def<sup>n</sup> for model parameters

- $t(f|e) \rightarrow$  cond<sup>n</sup>al. prob. of generating a French word  $f$  from an Eng. word  $e$ .
- $f \in F, e \in E \cup \{\text{NULL}\}$
- $q(j|i, l, m) \rightarrow$  prob. of alignment var.  $a_i$  taking the value  $j$ , conditioned on lengths  $l$  &  $m$ .

$$a_i \rightarrow j$$

$$i \in \{1, \dots, L\}, m \in \{1, \dots, M\}, j \in \{1, \dots, l\} \downarrow \\ (= \text{NULL})$$

Given these definitions:

$$p(f_1, \dots, f_m, a_1, \dots, a_n | e_1, \dots, e_l, i, m)$$

$$= \prod_{i=1}^m q(a_i | i, l, m) \cdot t(f_i | e_{a_i}) \quad (\text{v. Eng})$$

Fixes the index      Fixes the distribution of words

e.g.  $e =$  And the programme has been implemented

$f =$  le programme a été mis en application

Assumptions [V. Imp]

- ①  $q(a_i | i, l, m) \Rightarrow$   $(a_i)$  does not depend on any Eng. words or any other alignment values.
- ②  $t(f_i | e_{a_i}) \Rightarrow$  The word  $f_i$  is only dependent on the Eng. word it is aligned to

## Parameter Estimation in fully observed data

(10)

$(f^k, e^k, a^k) \leftarrow$  Training corpus consists of tuples

↳ alignment  
↳ English  
↳ French

$c(e, f) \rightarrow$  # times word  $e$  is aligned to word  $f$  in the training data

$c(e) \rightarrow$  # times  $e$  is aligned to any French word

$$t_{ML}(f|e) = \frac{c(e, f)}{c(e)}$$

ML: Maximum Likelihood

Remember: We are estimating this from the whole corpus  
(Not a single sentence)

$c(j|i, l, m) \rightarrow$  # times we see Eng. sent. of length  $l$ ,  
French sent. of len  $m$ , where word ' $i$ ' in  
Fr. sent. is aligned to word  $j$  in Eng

$$q_{ML}(j|i, l, m) = \frac{c(j|i, l, m)}{c(i, l, m)}$$

→ Algo 5.2 [from Notes]

## Applying IBM Model 2 [Imp.]

Once we have estimated  $t(f|e)$  &  $q(j|i,l,m)$

$$\Downarrow$$

$$p(f,a|e)$$

$$\arg \max_e p^{(e)} \stackrel{\Downarrow}{\leftarrow} p(f|e) = \sum_a p(f,a|e)$$

### Parameter Estimation

$$t(f|e) \quad q(j|i,l,m)$$

Data: parallel corpus

$$\begin{matrix} \text{pairs of sentences } f \rightarrow e \\ 500,000 \end{matrix} \quad f^1 \rightarrow e^1 \quad f^2 \rightarrow e^2 \quad \dots \quad f^N \rightarrow e^N$$

$$[f_1, f_2, \dots, f_m] \Leftrightarrow [e_1, e_2, \dots, e_l]$$

$$\uparrow \quad \quad \quad \uparrow \quad \quad \quad \uparrow \quad \quad \quad \uparrow$$

$$m \quad \quad \quad l$$

Issue: In the parallel corpus, we don't know the alignment for each training example.

$\Downarrow$   
Partially observable

Sol'n: ① HIRE annotators (Expensive)

## IBM Model - 1 : Alignments

- All alignments 'a' are equally likely :

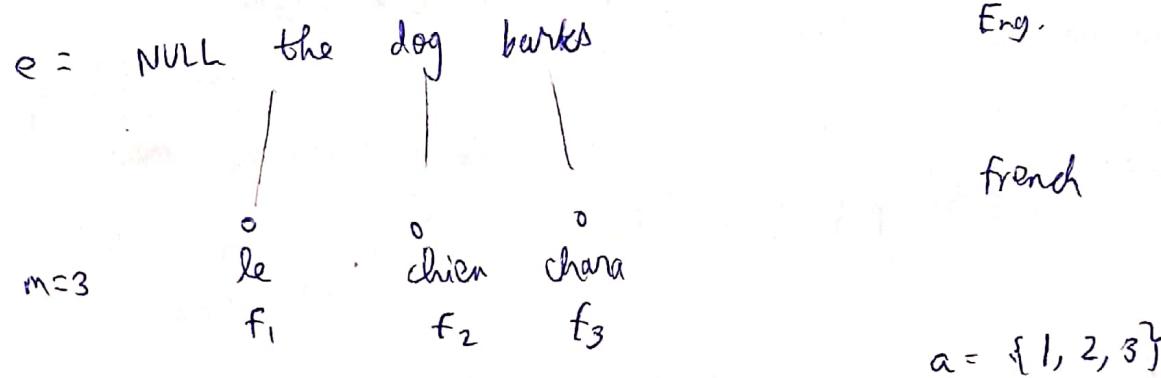
$$P(a|e, m) = \frac{1}{(l+1)^m}$$

### \* Translation probabilities

- Next step: come up with an estimate for :  $P(f|a, e, m)$
- In model 1, this is:

$$P(f|a, e, m) = \prod_{j=1}^m t(f_j | e_{aj})$$

e.g.



$$P(f|a, e, m) = t(\text{le} | \text{the}) \\ \times t(\text{chien} | \text{dog}) \\ \times t(\text{chara} | \text{barks})$$

### \* The Generative Process

To generate a French string f from an Eng. string e:

- Step-1: Pick an alignment 'a' with prob.  $\frac{1}{(l+1)^m}$

- Step-2: Pick the French words with prob.

$$P(f|a, e, m) = \prod_{j=1}^m t(f_j | e_{aj})$$

ie. The same French word may substitute for diff. Eng word in diff contexts

The final result:

$$P(f|e) P(f, a|e, m) = P(a|e, m) \times P(f|a, e, m) = \underbrace{\frac{1}{(l+1)^m}}_{\text{uni. prob. of alignment } a} \prod_{j=1}^m t(f_j | e_{aj})$$

Each fr. word an Eng. word is aligned to

e.g. How to estimate  $t(f|e)$ ?  $\rightarrow$  Through a parallel corpus

### \* IBM Model-2 (Extension of model-1)

- Only difference: We now introduce alignment or distortion parameters
- $q(a_j | j, l, m) = \text{Prob. that } j^{\text{th}} \text{ French word is connected to } i^{\text{th}}$  Eng. word, given send. length of e is l and f is m.
- Define:  $p(a | e, m) = \prod_{j=1}^m q(a_j | j, l, m)$   
where  $a = \{a_1, \dots, a_m\}$  is the alignment
- Result:  ~~$p(f | f)$~~   $p(F, a | e, m) = \prod_{j=1}^m q(a_j | j, l, m) \cdot t(f_j | e_{a_j})$

e.g.

$$l = 6$$

$$m = 7$$

$e = \text{And the prog. has been implemented}$  (len = 6)

$f = F_1 F_2 - F_3 \dots - F_7$  (len = 7)

$a = \{2, 3, 4, 5, 6, 6, 6\}$

$$p(a | e, 7) = q(2 | 1, 6, 7) \times q(3 | 2, 6, 7) \times q(4 | 3, 6, 7) \times q(5 | 4, 6, 7) \times q(6 | 5, 6, 7) \times q(6 | 6, 6, 7) \times q(6 | 7, 6, 7)$$

$$\begin{aligned} p(f | a, e, 7) &= t(F_1 | \text{the}) \times \\ &\quad t(F_2 | \text{program}) \times \\ &\quad t(F_3 | \text{has}) \times \\ &\quad t(F_4 | \text{been}) \times \\ &\quad t(F_5 | \text{implemented}) \times \\ &\quad t(F_6 | \text{implemented}) \times \\ &\quad t(F_7 | \text{implemented}) \end{aligned}$$

## The Generative Process

To generate a French string  $f$  from an Eng. string  $e$ :

- Step-1: Pick an alignment  $a = \{a_1, a_2, \dots, a_m\}$  with prob.:

$$\prod_{j=1}^m q(a_j | j, l, m)$$

- Step-2: Pick the French words with prob.

$$p(f | a, e, m) = \prod_{j=1}^m t(f_j | e_{a_j})$$

The final result:

$$p(f, a | e, m) = p(a | e, m) \cdot p(f | a, e, m) = \prod_{j=1}^m q(a_j | j, l, m) \cdot t(f_j | e_{a_j})$$

## Recovering Alignments

Given a sent. pair  $e_1, e_2, \dots, e_l$ ,  $f_1, f_2, \dots, f_m$  define

$$a_j = \arg \max_{a \in \{0, -1\}} q(a | j, l, m) \times t(f_j | e_a)$$

## \* EM Training of Models 1 and 2

### The Parameter Estimation Problem

Input:  $(e^{(k)}, f^{(k)})$  for  $k = 1 \dots n$

Eng. sent. Fr. sent.

O/P: parameters  $t(f | e)$  and  $q(i | j, l, m)$

challenge: We don't have alignments on our training examples  
(In ideal situation, sb would have annotated)

### Parameter Estm if the Alignments are Observed

$$t_{ML}(f | e) = \frac{\text{Count}(e, f)}{\text{Count}(e)}, \quad q_{ML}(j | i, l, m) = \frac{\text{Count}(j | i, l, m)}{\text{Count}(i, l, m)}$$

## Summary

- Key ideas in the IBM translation models:
  - Alignment variables
  - Translation parameters e.g.  $t(\text{chien} \mid \text{dog})$
  - Distortion parameters e.g.  $g(2 \mid 1, 6, 7)$
- The EM algorithm: an iterative algo for training the  $g$  and  $t$  parameters.
- Once the parameters are trained, we can recover the most likely alignments on our training examples.

\* Recommender Systems: reduce info. overload by estimating relevance

- Paradigms:
- ① Personalized recommendations
  - ② Collaborative (popular among peers) ✓
  - ③ Content-based: (More of what I've liked)
  - ④ Knowledge-based: (what fits based on my needs)

Hybrid: Comb" of various inputs and/or composition of diff. mechanisms

\* Collaborative Filtering (CF) → pros: No knowledge engg' reqd.

Approach: Use "wisdom of the crowd" to recommend items

Basic assumptions (v. imp.) (2)

- ① Users give ratings to catalog items (implicitly or explicitly)
- ② Customers who had similar tastes in the past, will have similar tastes in the future.

Pure CF approaches:

I/p: Only a matrix of given user-item ratings

o/p:

- A (numerical) pred" indicating to what degree the current user will like or dislike certain items
- A top-N list of recommended items

① User-based NN-CF

• "Active user" → Alice      Item+J (Not yet seen by Alice)

pred": Alice's rating for item J:

= Avg. of the ratings of her peers who liked the same items as Alice in the past and who have rated item J.

- Do this for all items not seen by Alice and recommend the best-rated.

### (a) Measuring user similarity

#### Pearson correlation (g) →

• Linear correlation

• Lies b/w -1 and 1  
+ neg.  
↳ Total +ve linear corr.

a, b: Users

$r_{a,p}$ : Rating of user a for item p

p: Set of items, rated both by a and b

$$g = \text{sim}(a, b) = \frac{\sum_{p \in P} (r_{a,p} - \bar{r}_a)(r_{b,p} - \bar{r}_b)}{\sqrt{\sum_{p \in P} (r_{a,p} - \bar{r}_a)^2} \cdot \sqrt{\sum_{p \in P} (r_{b,p} - \bar{r}_b)^2}}$$

0  
↳ No corr.

★ Note: g is calculated on raw user ratings

Reason: Users have different bias while rating.

Some rarely give 5\* while some do it easily (frequently).

→ measures relative distance from the avg. rating

Numerator: COVARIANCE

Denominator: Std. Dev. of user A \* Std. Dev. of User B

#### Pearson correlation

- Takes differences in rating behv. into account
- works well in usual domains, compared with alternate measures such as cosine similarity.

### (b) Making Predictions

$$\text{pred}(a, p) = \bar{r}_a +$$

↓  
User's avg.  
rating

$$\frac{\sum_{b \in N} \text{sim}(a, b) * (r_{b,p} - \bar{r}_b)}{\sum_{b \in N} \text{sim}(a, b)}$$

Neighbour's bias  
Weighted similarity  
Normalization  
Neighbours decided on similarity threshold

### IMPROVING PREDICTION

Better similarity, weighting matrices and neighbour selection

- Give more weights to items that have a higher variance.
- # Co-rated items : Use "significance weighting"  
e.g. linearly reducing the weight when the # co-rated items is low.
- Case amplification : Give more weight to "very similar" neighbours i.e. where similarity value is close to 1.
- Nbd. selection : Use similarity threshold or fixed No. of neighbours.

### \* ITEM BASED COLLABORATIVE FILTERING

Basic idea: Use similarity b/w items (and not users) to make predictions

Cosine similarity → produces better results in item-to-item filtering

Ratings → n-dimensional space

$$\text{sim}(\vec{a}, \vec{b}) = \frac{\vec{a} \cdot \vec{b}}{|\vec{a}| \cdot |\vec{b}|}$$

Issue (v. imp.) : Does not take into account user's bias.

## Adjusted cosine similarity

Avg. user rating taken into account

$$\text{sim}(a, b) = \frac{\sum_{u \in U} (r_{u,a} - \bar{r}_u)(r_{u,b} - \bar{r}_u)}{\sqrt{\sum_{u \in U} (r_{u,a} - \bar{r}_u)^2} \sqrt{\sum_{u \in U} (r_{u,b} - \bar{r}_u)^2}}$$

## Making predictions

$$\text{pred}(u, p) = \frac{\sum_{i \in \text{ratedItems}(u)} \text{sim}(i, p) * g_{u,i}}{\sum_{i \in \text{ratedItems}(u)} \text{sim}(i, p)}$$

Most real-life datasets: 20-50 neighbours work well.

## Pre-processing for item-based filtering [That's why they are preferred]

Approach:

- calculate all pairwise item similarities in advance.
- The nbd. to be used at run-time is typically small
- Item similarities are supposed to be more stable than user similarities.

Memory Req.:

( $N \approx 1$  million)

- # items =  $N$
- $\Rightarrow N^2$  pairwise similarities to be memorized

$$\begin{matrix} I_1 & I_2 & I_3 & \dots \\ T_1 & & & \\ T_2 & & & \\ T_3 & & & \\ \vdots & & & \end{matrix}$$

- Easily expandable ✓

In practice, matrix would be sparse, so space can be further reduced.

stable  
(won't change much)

while user based matrix changes on small change

## ④ ISSUES WITH CF (collaborative filtering)

### \* Data Sparsity Problems

- Cold start problem

- How to recommend new items? What to recommend to new users?

(soln)

- o straightforward approach: (Imp.)

- Force users to rate a set of items.

{  
in  
initial  
phase}

- Use other methods like content-based, demographic or non-personalized in the initial phase.

### The Cold-start Problem [Coursera] - Univ. of Minnesota

- New user
- New item
- New system

#### e. New user

- problem: No profile, No preferences

- No issue if system is non-personalized

- Else: provide useful default personalized options

e.g.

- o Offer popular items — Get data from its ratings
- o Demographically relevant (or age, sex) — Domain dependent

- o product association — Browsing history or shopping cart
- o Trust-based social network data (or Referral data)

- New items
    - Challenge:
      - can't recommend - because we don't know whether user will like or not
      - Need ratings quickly
  - Options
    - use content-based approach (Degrading our rec. system)
    - Recommend to random, (or) well-chosen set of users  
(people with diverse tastes, tolerance, interest, influence)
  - New systems
    - syndicated data (get data from another source)
      - ↳ could be content data or rating data
    - Design a non-recommender system that captures data to feed a later recommender
- [End-of-Coursera]

- \* Alternative approach to Data sparsity problems
    - "Transitivity" of neighbourhoods
    - ② Recursive CF: Assume there is a v-close neighbour  $n$  of  $u$  who has not rated the item  $i$  yet  
(This approach works better when data is sparse)
    - Graph-based methods
      - \* Spreading activation: Exploit the supposed "transitivity" of customer tastes and thereby augment the matrix with additional info.
      - Idea: Use paths of length  $> 3$  (Extensible) to recommend items.
- [Nearest nb. based CF]

## Association Rule Mining

- Market-basket ( $\text{Baby food} \Rightarrow \text{Diaper(M)}$ )

### Algo:

Trans.  $t_i = \{\text{set of items purchased during that transaction}\}$

- Detect rules ( $A \Rightarrow B$ ) from  $D = \{t_1, t_2, \dots, t_n\}$

- Two qty: **support** and **confidence** (Both should be high)

- Let  $X$ : set of items

$c(X)$ : How many trans. ( $t_i$ ),  $X$  is present

$$\text{Support} = \frac{c(X \cup Y)}{|D|}$$

$$\text{Confidence} = \frac{c(X \cup Y)}{c(X)}$$

Total # trans.

- ← Simple approach: Transform 5-point ratings into binary ratings  
(1 = above user avg.)

### Making recommendations

USER's

- Determine "relevant" rules based on transactions

- Determine items not already bought by user.

- Sort the items based on the rules' confidence values.

## \* Evaluating Recommender Systems

### (3 types of Evaluation Metrics)

#### ① Comparing Values

- measure how close the predicted ratings are to the true user ratings (for all the ratings in the test set)
- Mean Absolute Error (MAE)

$$MAE = \frac{1}{n} \sum_{i=1}^n |p_i - r_i|$$

predicted rating  
actual rating

- Root Mean Square Error (RMSE) : similar to MAE, but places more emphasis on larger deviation

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (p_i - r_i)^2}$$

#### ② Comparing Recommendations : precision and recall

$$\text{precision} = \frac{tp}{tp + fp}$$

$$\text{recall} = \frac{tp}{tp + fn}$$

Rec.  $\approx$  Info. Retr. task

$$(h) \text{ prec} = \frac{|\text{good movies rec.}|}{|\text{all rec.}|}$$

$$\text{recall} = \frac{|\text{good movies rec.}|}{|\text{all good movies}|}$$

Typically, a rec. system is tuned to increase precision  
(Recall decreases as a result)

#### ③ Metric<sup>s</sup> : rank-pos<sup>n</sup> matters

- Rank metrics extend recall and precision to take the pos<sup>n</sup> of correct items in a ranked list into account.

- o Relevant items are more useful when they appear earlier in the rec. list
- o particularly imp. in rec. systems as lower ranked items may be overlooked by users

\* Metric: Normalized Discounted Cumulative Gain

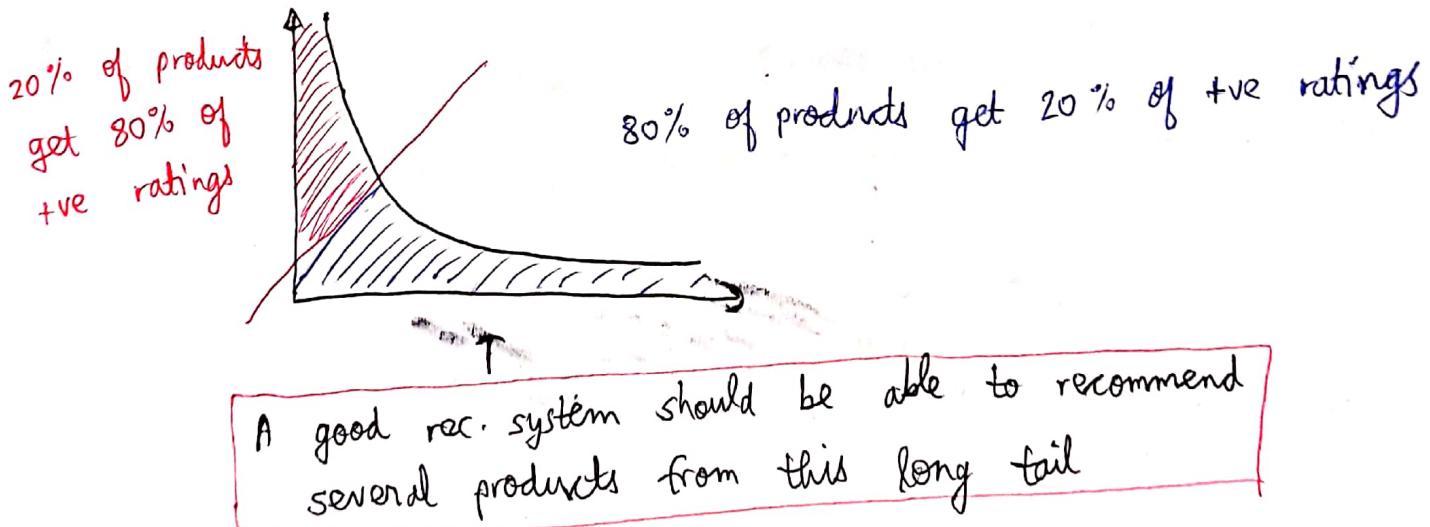
$$DCG_{pos} = \text{rel}_1 + \sum_{i=2}^{pos} \frac{\text{rel}_i}{\log_2 i} \quad (\text{Logarithmic reduction factor})$$

nDCG

- Divide by the ideal recommendations DCG
- Normalize to the interval [0..1]

(Normalize DCG at rank n by the DCG value at rank n of the ideal ranking)

\* An Imperfect World : Concentration Effects (long tail)



24/04/2020

## Multimedia IR

Appl.: Medical Images  
(X-rays)

### \* Use cases:

- Images
  - Find look-alike pictures (Google Reverse Img search) ↗
  - Recognize landmarks (Google Lens)
  - People wearing white suits (text-based img. retrieval) → Most common
- Two categories :
  - ① Img to Img search
  - ② Text to Img. search
- Audio
  - Finding similar sounds [Music to Music]
    - Music
    - Animal sounds
    - copyrights
  - Retrieve relevant speech or podcasts or music [Text to Music]
- Videos
  - Finding similar videos [Video to Video]
    - Copyrights (e.g. if someone uploads full movie on YT)
  - Retrieve relevant news telecasts
  - Find specific actions: Kathak dance } [Text to video]
  - Lecture

### \* Basic IR Working Intuition

Doc: [Images]

q: [Images]

Case-1      2 alike  
Doc & q → both are multimedia

Doc → Rep.  
q → Rep.

This representation should be comparable through some similarity measure

Doc: [Images]

q: [Text]

Img  $\rightarrow$  RGB  $\rightarrow$   $N \times M \times 3$

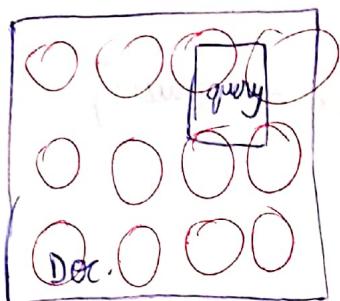
Even if  $N = M = 100$ ,  $N \times M \times 3 = 30k$  is huge # parameters to compare.

What we do? : Img is treated as a tensor of  $N \times M \times 3$

↓  
Encoded into lower-dimensional vector (100)

Now, vector-based comparison b/w doc & query.

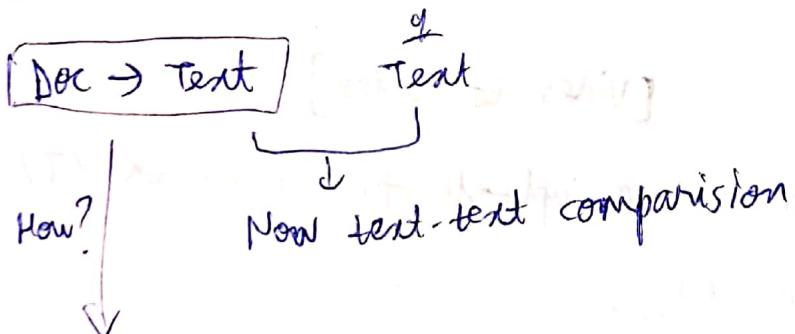
Challenge:



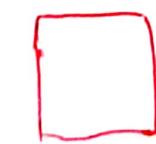
→ query is a small part of doc.

Divide img (doc) into smaller frags and then encode

Case-II: Doc is multimedia (Img) and query  $\neq$  Text.



Suppose: Image www (Source  $\rightarrow$  Webpages)



Caption/title

(while crawling & retrieve it)

Image Captioning

Android photo



ML classifier



Assign tag to image

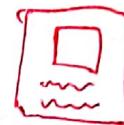
IR system now uses this tag.

→ Same logic can be applied to audio and video.

### \* Piggyback Retrieval [Imp.]

- Convert everything to text
- Use text based retrieval

TITLE  
↑



- Use metadata or loose annotations (automated or manual)
  - Image caption (text surrounding images on web page)

### \* Info. Loss (while converting img/audio/video → text)

Because it may not be annotated considering all possible scenarios.

#### • Semantic gap

Features: segmented-blobs; Salient regions, pixel-level histograms, Fourier descriptors etc.

#### \* Example: Sub-image matching

## 27/07/2020 CONTENT BASED RETRIEVAL (Multimedia IR)

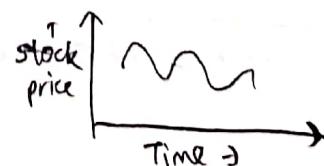
Objective: Design a fast searching method that will search database of multimedia objects to locate objects that matches the query object, exactly or approximately.

e.g.

- Img. similar to a given query img.
- companies whose stock prices move similarly

$y_{obj}^{MM}$ : time series data

[same stock price variation]



- X-ray img. that contains sth. that has a texture of a tumour.

## TERMINOLOGIES

QUESTION

→ Distance or Dissimilarity

$$D(O_1, O_2) = \text{Dist. b/w 2 obj. } O_1 \text{ & } O_2$$

e.g. Euclidean dist.

→ Whole match: Given a collection of N obj:  $O_1, O_2 \dots O_n$  and a query Q, find those obj. that are within dist. E from Q.

Note: Query and objs are of same type : img-img  
audio-audio etc.

→ Sub-pattern match : Query is allowed to specify only part of the object.

- e.g. • Obj's. are  $512 \times 512$  or larger size img. (medical X-rays)
- Query is  $32 \times 32$  img. (X-ray of a tumor)

## \* 4 Key Requirements (V. Imp)

- Fast (distance calculation)
- Correct (Return all qualifying objects. — No false negatives)
  - Note: False positives — acceptable  
They can be discarded in the post-processing step
- Small space overhead
- Dynamic (Easy to insert, delete and update objects)

## \* GEMINI (GENERIC Multimedia object INDEXING)

Two key ideas: (To avoid 2 disadvs. of sequential scanning)

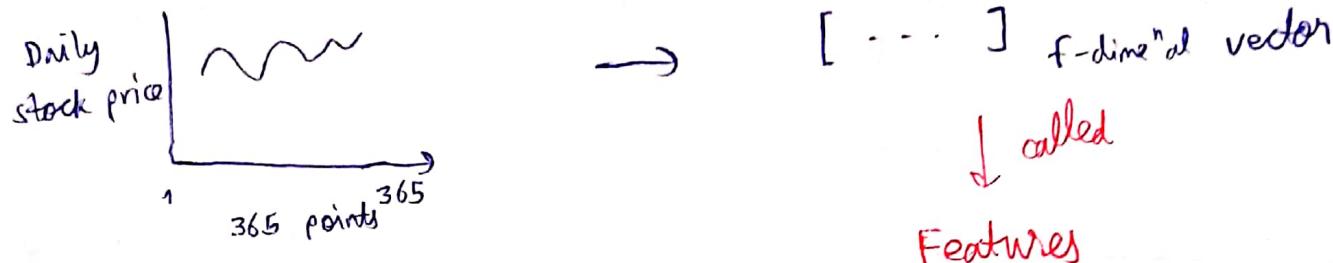
- 1) a 'quick-and-dirty' test, to discard quickly the vast majority of non-qualifying objects (possibly allowing some false alarms)
- 2) use of spatial access methods, to achieve faster-than-sequential searching.

## Quick and Dirty test (QD test) - Idea

To categorize obj. with a single or multiple numbers, which can help us discard many non-qualifying obj. sequences.

e.g.

If obj. is a series of Nos., a single number can be avg.



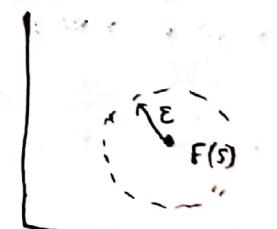
(Numbers that contain some info. about the multimedia object)

Def<sup>n</sup> :  $F(\cdot) \rightarrow$  mapping of obj. to f-dim. pt.  
 $F(Q) \rightarrow$  f-D pt.

Then use feature space to apply QD test.

### \* Algo 1 : Search

1. Map the query obj.  $Q$  into a pt.  $F(Q)$  in feature space.
2. Using a spatial access method, retrieve all points within a desired tolerance  $\epsilon$  from  $F(Q)$ .
3. Retrieve the corresponding objects, compute their actual distance from  $Q$  and discard the false alarms.  
 (By computing actual distance  $D$  of these points)



### \* No False Negatives Guarantee

Lemma (Lower Bounding) : To guarantee no false dismissals for whole-match queries, the feature extraction fn  $F(\cdot)$  should satisfy

the foll. formula:

$$D_{\text{feature}}(F(O_1), F(O_2)) \leq D(O_1, O_2)$$

## \* Spatial Access Methods (SAM)

The mapping provides the key to improve the 2<sup>nd</sup> drawback of sequential scanning.

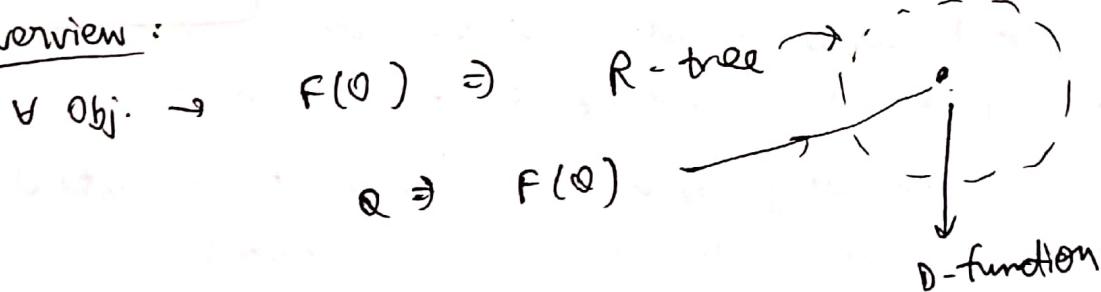
↳ R-trees (A data-structure)

↓  
an index tree-structure derived from the B-tree that uses  
multi-dimensional indexes

## \* Algo 2: (GEMINI)

1. Determine the distance  $f^n D()$  b/w 2 objects
2. Find 1 or more numerical feature-extraction functions, to provide a 'quick-and-dirty' test.
3. Prove that the distance in feature space lower-bounds the actual distance  $D()$ , to guarantee correctness.
4. Use a SAM (e.g. an R-tree) to store and retrieve the F-D feature vectors.

Overview:



## \* More about D() and F()

D(): Distance function  
F(): Feature function

### "Good" dist. $f^n$

- A domain expert decides (independent of GEMINI)
- The GEMINI methodology focuses on speed of search only.

### "Good" feature $f^n$

- should facilitate step-3 (distance lower-bounding)
- should capture most of the ~~et~~ characteristics of the objects.

## \* $g = D$ Time Series (e.g.)

- Variation of stock prices
- Dist.  $f^n$ : Euclidean dist., Time warping

→ Bad feature: Only 1<sup>st</sup> day value of the company stock

→ Good feature:  

- Avg.
- coeff. of DFT (Discrete Fourier Trans.  
gives n-coeff., keep top-k coeff.)

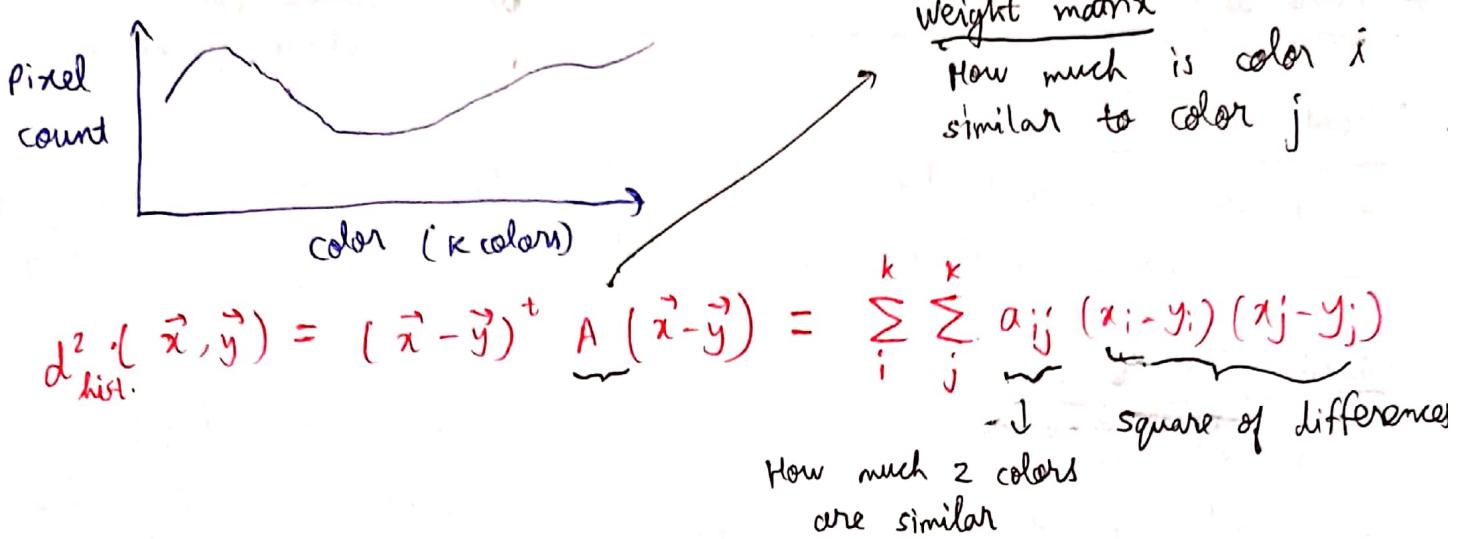
$$F_{\text{DFT}}^{(0)} =$$

Parseval's theorem: coeff. of DFT and dist. computed in feature space using these coeff. will be a lower bound on the Euclidean distance.  
(Hence GEMINI algo can be used)

## \* 2-D Color Images

Distance: k-element color histogram

[Take 64 or 32 diff. colors. Given an img., plot a histogram]  
x-axis = colors, y-axis: # pixels



### 3-D feature vector

$$\vec{x} = (R_{\text{avg.}}, G_{\text{avg.}}, B_{\text{avg.}})^T$$

$$G_{\text{avg.}} = \left( \frac{1}{P} \right) \sum_{p=1}^P G(p) \quad | \quad R_{\text{avg.}} = \left( \frac{1}{P} \right) \sum_{p=1}^P R(p) \quad | \quad B_{\text{avg.}} = \left( \frac{1}{P} \right) \sum_{p=1}^P B(p)$$

$$d_{\text{avg.}}^2(\vec{x}, \vec{y}) = (\vec{x} - \vec{y})^T (\vec{x} - \vec{y})$$

Lower Bound : Quadratic Distance Bounding Theorem

### Conclusion

Two key elements to speed up the retrieval of multimedia objects are:

- QD test       $\leftarrow F()$
- Spatial Access methods (SAM)       $\leftarrow R\text{-Trees}$