Proposal Marks are Out!
Assign 1, Individual, Tomorrow 7 AM EST

**Blackboard will return to service at 6 am EST.**

Blackboard is unavailable every weekday (Monday to Friday) from 5 am to 6 am EST for regular maintenance.
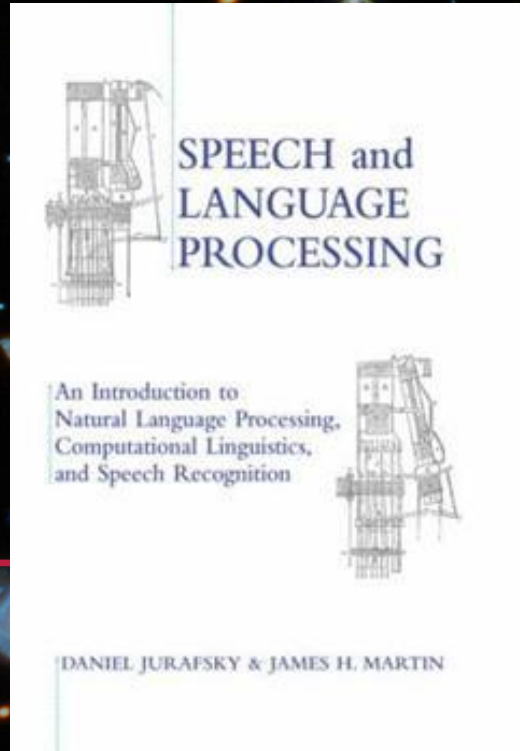
At 6 am you can reload this page to access Blackboard.

Office Hour
2:30 PM – 3:30 PM
Lambton Tower 5111

*n*-Gram Language Models

# Language Modeling

CH04

Language Model
$Context_i \longrightarrow Context_{i+1}$

# *n*-Gram Language Model

## Context Window of Size *n*

Recent Past of Size *n-1* → Future of Size 1

$$W_{i+1} \dots W_{i+n-2} \; W_{i+n-1} \longrightarrow W_{i+n}$$

# *n*-Gram Language Model

$$\longrightarrow W_{i+1}$$ 1-gram = unigram

$$W_{i+1} \longrightarrow W_{i+2}$$ 2-gram = bigram

$$W_{i+1}\, W_{i+2} \longrightarrow W_{i+3}$$ 3-gram = trigram

$$W_{i+1} \dots W_{i+n-2}\, W_{i+n-1} \longrightarrow W_{i+n}$$ n-gram

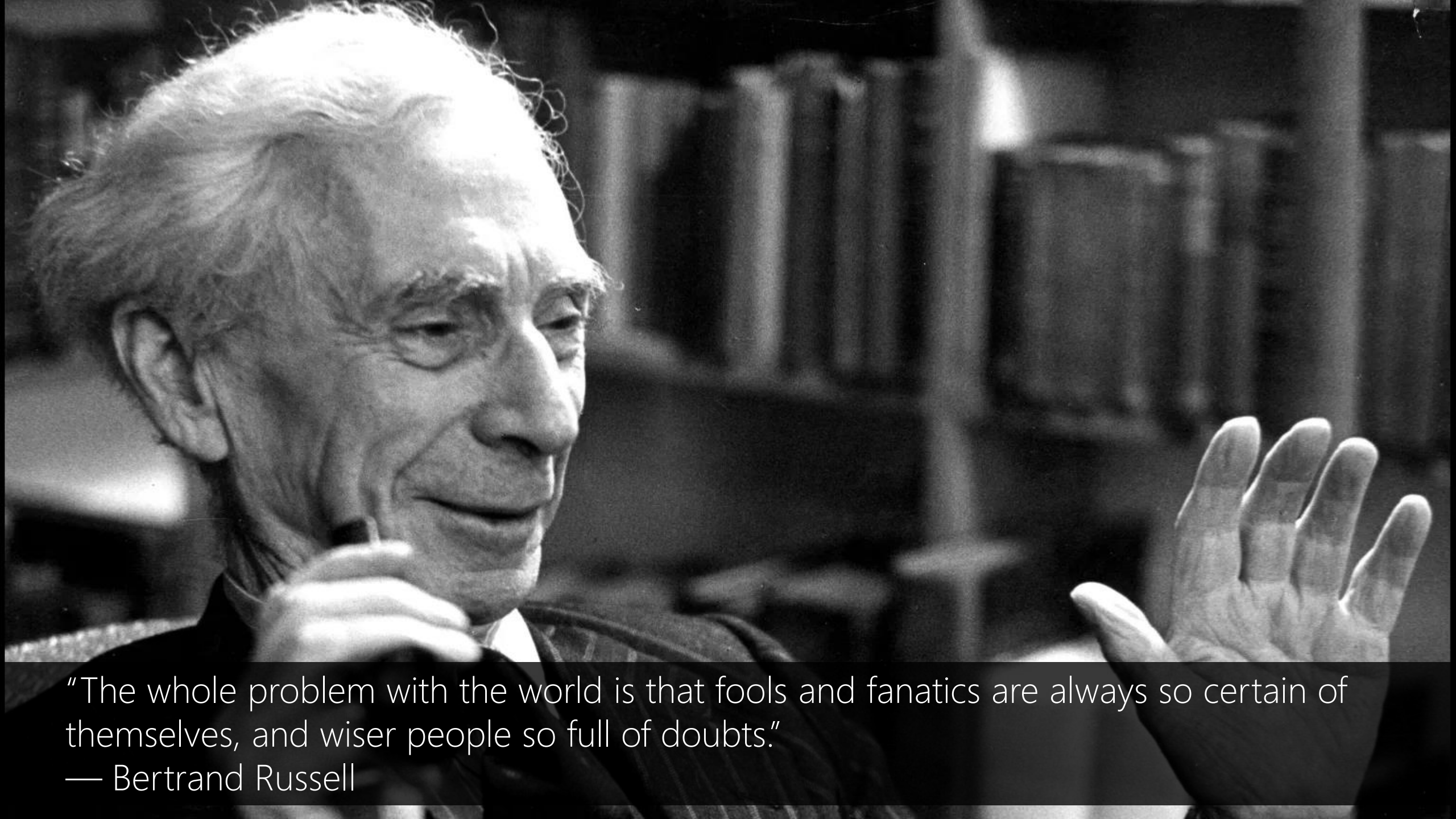# Frequentist Probability
*as opposed to Bayesian Probability*

---

*Frequentist probability or frequentism is an interpretation of probability that defines an event's probability as the limit of its relative frequency in many trials* - Wikipedia

# *n*-Gram Language Modeling

Recent Past of Size *n-1* → Future of Size 1 → Most Frequent Future Given the Past

$$w_{i+1} \dots w_{i+n-2} \, w_{i+n-1} \rightarrow w_{i+n} = \text{Max } P(w \mid w_{i+1} \dots w_{i+n-2} \, w_{i+n-1}) \text{ in all } w \in V$$

$$P(w \mid w_{i+1} \dots w_{i+n-2} \, w_{i+n-1}) = \frac{P(wi_{+1} \dots w_{i+n-2} \, w_{i+n-1} w)}{P(wi_{+1} \dots w_{i+n-2} \, w_{i+n-1})}$$

$$= \frac{\#(wi_{+1} \dots w_{i+n-2} \, w_{i+n-1} w)}{\#(wi_{+1} \dots w_{i+n-2} \, w_{i+n-1})}$$

"The whole problem with the world is that fools and fanatics are always so certain of themselves, and wiser people so full of doubts."
— Bertrand Russell

# Chain Rule of Probability

$$P(w_1 w_2 \dots w_n) = P(w_1)\, P(w_2 \mid w_1)\, P(w_3 \mid w_1 w_2) \dots P(w_n \mid w_1 w_2 w_3 \dots w_{n-1})$$

$$= \prod_{k=1}^{n} P(w_k \mid w_1 \dots w_{k-1})$$

$$= \prod_{k=1}^{n} P(w_k \mid w_1^{\,k-1})$$

# *Approximation* to Chain Rule

## Generalizability

Language is creative! A particular context might have never occurred before!

# *Approximation* to Chain Rule

## Efficiency

probability of a word given entire history, approximate the history by just the last few words

# Unigram Approx.

*Bag-of-Word (BoW). Why?*

$$P(w_1 \, w_2 \, ... \, w_n) = P(w_1)P(w_2|w_1)P(w_3|w_1w_2) \, ... \, P(w_n|w_1w_2w_3...w_{n-1})$$

$$= P(w_1)P(w_2| \quad )P(w_3| \quad ) \, ... \, P(w_n| \quad )$$

$$= P(w_1)P(w_2)P(w_3) \, ... \, P(w_n)$$

# Bigram Approx.

*Markovian: probability of a variable depends only on the previous variable*

$$P(w_1 \, w_2 \, ... \, w_n) = P(w_1)P(w_2|w_1)P(w_3|w_1 w_2) \, ... \, P(w_n|w_1 w_2 w_3 ... w_{n-1})$$

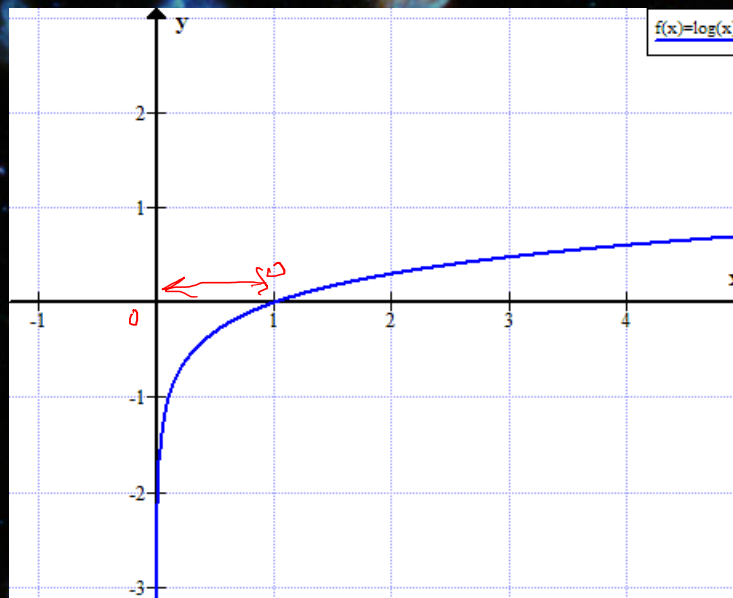$$= P(w_1)P(w_2|w_1)P(w_3| \quad w_2) \, ... \, P(w_n| \qquad \qquad w_{n-1})$$

# Trigram Approx.

$$P(w_1\ w_2 \ldots w_n) = P(w_1)P(w_2|w_1)P(w_3|w_1w_2) \ldots P(w_n|w_1w_2w_3\ldots w_{n-1})$$
$$= P(w_1)P(w_2|w_1)P(w_3|w_1w_2) \ldots P(w_n|\qquad w_{n-2}w_{n-1})$$

# Approx. n-gram Language Modeling

## Corpus: Brown University

['The', 'Fulton', 'County', 'Grand', 'Jury', 'said', 'Friday', 'an', 'investigation', 'of', ..., '.'],
['The', 'jury', 'further', 'said', 'in', 'term-end', 'presentments', 'that', 'the', 'City', ... 'conducted', '.'],
['The', 'September-October', 'term', 'jury', 'had', 'been', 'charged', 'by', 'Fulton', 'S...', 'Allen', 'Jr.', '.'],
['``', 'Only', 'a', 'relative', 'handful', 'of', 'such', 'reports', 'was', 'received', "''", ',', '... 'city', "''", '.'],
['The', 'jury', 'said', 'it', 'did', 'find', 'that', 'many', 'of', "Georgia's", 'registration', ... 'ambiguous', "''", '.']

| Make it worse! | Gives chance to new combination |
|---|---|
| P([Mr.][and][Mrs.]) | P([Mr.][and][I]) |
| 0.00045851027827207127 | 0.0 |
| 1.42083315097917766e-05 | 1.75171210394693e-06 |
| 9.078228423943108e-08 | 6.422936315754214e-08 |

# Log of Probabilities

$P(x_1) \times P(x_2) \times \ldots \times P(x_n) \propto \log P(x_1) + \log P(x_2) + \ldots + \log P(x_n)$

left and right sides have same order!

$[0, 1] \longrightarrow [-\infty, 0]$

Product $\longrightarrow$ Sum

# Self-supervised

Self-supervised learning is the key to AI understanding the world

Yann LeCun: Dark Matter of Intelligence and Self-Supervised Learning | Lex Fridman Podcast #258
https://www.youtube.com/watch?v=SGzMElJ11Cc

Tenet, Christopher Nolan, 2020
Budget $200 m
Box office $363 million

# *n*-Gram Language Modeling

Recent Past → Current ← Recent Future

---

$$W_{i+1} \ldots W_{i+n-2} \; W_{i+n-1} \longrightarrow W_{i+n} \longleftarrow W_{i+n+1} \; W_{i+n+2} \ldots W_{i+n+j}$$

Following, Christopher Nolan (1998)
Budget $6,000
Box office $48,482

*Evaluating* Language Models

*Evaluating* Language Models

Higher *n* in n-gram, the better?
More history, the better prediction of future?

*Evaluating* Language Models

Qualitative → Let's Communicate → Generate

| | |
|---|---|
| **1** gram | –To him swallowed confess hear both. Which. Of save on trail for are ay device and rote life have <br><br> –Hill he late speaks; or! a more to leg less first you enter |
| **2** gram | –Why dost stand forth thy canopy, forsooth; he is this palpable hit the King Henry. Live king. Follow. <br><br> –What means, sir. I confess she? then all sorts, he is trim, captain. |
| **3** gram | –Fly, and will rid me these news of price. Therefore the sadness of parting, as they say, 'tis done. <br><br> –This shall forbid it should be branded, if renown made it empty. |
| **4** gram | –King Henry. What! I will go seek the traitor Gloucester. Exeunt some of the watch. A great banquet serv'd in; <br><br> –It cannot be but so. |

**Figure 3.3** Eight sentences randomly generated from four n-grams computed from Shakespeare's works. All characters were mapped to lower-case and punctuation marks were treated as words. Output is hand-corrected for capitalization to improve readability.

# Cross Evaluating Language Models

Biased toward the corpus! Dialect, Genre, …

Better LM is the one that can generalize!

Models

M* $P(L \mid M*) = 1$

$P(l \mid M*) = ?$

$B.S.$

M_2 $P(L \mid M_2) = 0.001$

$P(l \mid M_2) = 1$

M_1 $P(L \mid M_1) = 0$

$P(l \mid M_1) = 0$

POOR

Language L

$l$

# Find M*
## Golden Model

Assumptions:
- The language space is available. ✗
- Search the model space to find M*, assuming it exists!

# Find M*

## Golden Model

Relax Assumptions:
- The language space is available. → Random Subsets
- Search the model space to find M*, assuming it exists!

Language L

*l*

Model M*
P(L | M*) = 1
P(*l* | M*) = 1

# Find ~~M*~~ M^

~~Golden~~ Silver Model

Relax Assumptions:
- The language space is available. → Random Subsets
- Search the model subspace to find M^, assuming it exists!

$P(L \mid M*) = 1$

$P(l \mid M*) = 1$

Models

M*

Language L

$l$

$M_1$

$P(l \mid M_1) = 0$

$P(L \mid M_1) = 0$

Models

$P(L \mid M^*) = 1$

$P(\ell \mid M^*) = 1$

M*

Language L

$\omega_1$

$\ell$

$\omega_1$

$P(\ell \mid M_4) = 0.9$

$P(L \mid M_4) = ?$

n-Gram Models

5-gram

$M_4$

$M_1$

$P(\ell \mid M_1) = 0$

$P(L \mid M_1) = 0$

$\omega_1 \omega_2$

$0.9$

$0.9 \cdot 0.9$

$\ell \subset L$

$$\hat{M} = \underset{M \in \text{Models}}{\arg\max} \, P(\ell \mid M)$$
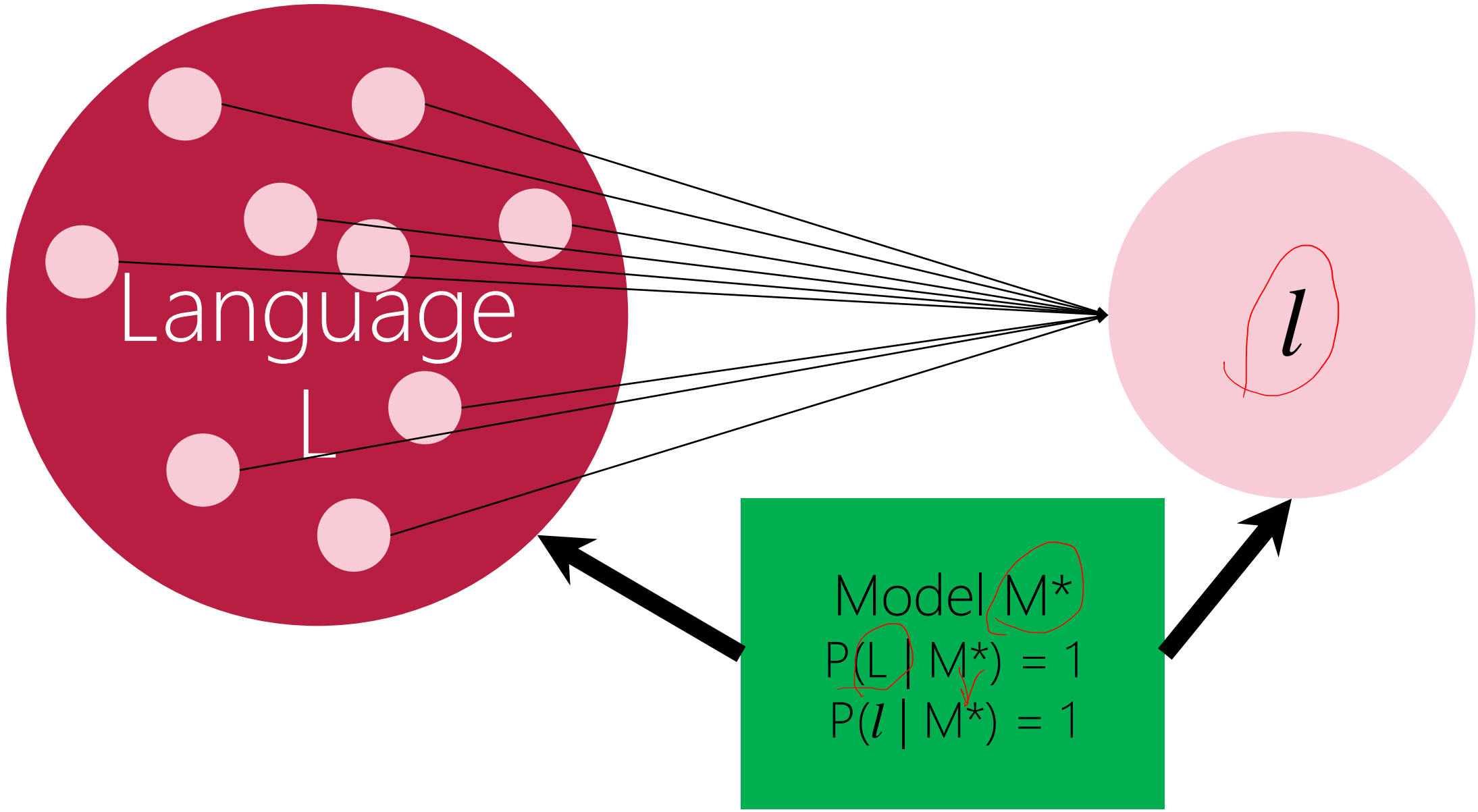
Relax Assumptions:
- The language space is available. → Random Subsets
- Search the model subspace to find $\hat{M}$, assuming it exists!

Likelihood

$$\hat{M} = \underset{M \in \text{Models}}{\text{argmax}} \; P(l \mid M)$$

Maximum Likelihood Estimation (MLE)

$$\hat{M} = \underset{M \in \text{Models}}{\text{argmax}} \underbrace{\mathcal{L}(l \mid M)}_{\text{Likelihood}}$$

Likelihood

Maximum Likelihood Estimation (MLE)

$$\hat{M} = \underset{M \in \text{Models}}{\operatorname{argmax}} \frac{P(l, M)}{P(M)}$$

$P(l/M)$

$$\hat{M} = \underset{M \in \text{Models}}{\text{argmax}} \frac{P(l, M)}{P(M)}$$

$P(M) \sim$ Uniform Distribution (equal chance)

$$\hat{M} = \underset{M \in \text{Models}}{\arg\max} P(l, M)$$

$$P(l, M) = P_M(l) = \mathcal{L}_M(l)$$

$$\hat{M} = \underset{M \in \text{Models}}{\arg\max} \mathcal{L}_M(l)$$

# Likelihood for a Language Model

$l$: ] The ][ course ][ COMP8730 ][ is ][ about ][ nlp ][ . ][ The ][ instructor ][ 's ][ name ][ is ][ Hossein ][ . ][ There ][ are ][ 13 ][ students ][ in ][ the ][ class ][ . ]

M1 = |token|-gram model = 22-gram model → $P_{22\text{-gram}}(l) = \mathcal{L}_{22\text{-gram}}(l) = 1$

We found our silver model!
M^ = M1

# Likelihood for a Language Model

$l$: [ The ][ course ][ COMP8730 ][ is ][ about ][ nlp ][ . ][ The ][ instructor ][ 's ][ name ][ is ][ Hossein ][ . ][ There ][ are ][ 13 ][ students ][ in ][ the ][ class ][ . ]

M1 = |token|-gram model = 22-gram model → $P_{\text{22-gram}}(l) = \mathcal{L}_{\text{22-gram}}(l) = 1$

M2 = |vocab|-gram model = 16-gram model → $P_{\text{16-gram}}(l) = \mathcal{L}_{\text{16-gram}}(l) = ?$

$l$: [ The ][ course ][ COMP8730 ][ is ][ about ][ nlp ][ . ][ The ][ instructor ][ 's ][ name ][ is ][ Hossein ][ . ][ There ][ are ][ 13 ][ students ][ in ][ the ][ class ][ . ]

M1 = |token|-gram model = 22-gram model $\rightarrow$ P$_{22\text{-gram}}(l)$ = $\mathcal{L}_{22\text{-gram}}(l)$ = 1

M2 = |vocab|-gram model = 16-gram model $\rightarrow$ P$_{16\text{-gram}}(l)$ = $\mathcal{L}_{16\text{-gram}}(l)$ = ?

P( [ The ][ course ][ COMP8730 ][ is ][ about ][ nlp ][ . ][ The ][ instructor ][ 's ][ name ][ is ][ Hossein ][ . ][ There ][ are ][ 13 ][ students ][ in ][ the ][ class ][ . ] )

$= P(the) ($

$= P(course | the) ($  $| the) ($  $(the, cou...)$

# Chain Rule of Probability

$P(w_1 w_2 \ldots w_n) = P(w_1) \, P(w_2 | w_1) \, P(w_3 | w_1 w_2) \ldots P(w_n | w_1 w_2 w_3 \ldots w_{n-1})$

$$= \prod_{k=1}^{n} P(w_k | w_1 \ldots w_{k-1})$$

$$= \prod_{k=1}^{n} P(w_k | w_1^{k-1})$$

# Likelihood for a Language Model

$l$: [ The ][ course ][ COMP8730 ][ is ][ about ][ nlp ][ . ][ The ][ instructor ][ 's ][ name ][ is ][ Hossein ][ . ][ There ][ are ][ 13 ][ students ][ in ][ the ][ class ][ . ]

M1 = |token|-gram model = 22-gram model $\rightarrow$ P$_{22\text{-gram}}(l) = \mathcal{L}_{22\text{-gram}}(l) = 1$

M2 = |vocab|-gram model = 16-gram model $\rightarrow$ P$_{16\text{-gram}}(l) = \mathcal{L}_{16\text{-gram}}(l) = ?$

P$_{\text{([ The ][ course ][ COMP8730 ][ is ][ about ][ nlp ][ . ][ The ][ instructor ][ 's ][ name ][ is ][ Hossein ][ . ][ There ][ are ][ 13 ][ students ][ in ][ the ][ class ][ . ])} =$
P$_{\text{([ The ])}}$P$_{\text{([ course ][ COMP8730 ][ is ][ about ][ nlp ][ . ][ The ][ instructor ][ 's ][ name ][ is ][ Hossein ][ . ][ There ][ are ][ 13 ] …. |[ The ])}}$

# Likelihood for a Language Model

$l$: [ The ][ course ][ COMP8730 ][ is ][ about ][ nlp ][ . ][ The ][ instructor ][ 's ][ name ][ is ][ Hossein ][ . ][ There ][ are ][ 13 ][ students ][ in ][ the ][ class ][ . ]

M1 = |token|-gram model = 22-gram model $\rightarrow$ $P_{\text{22-gram}}(l) = \mathcal{L}_{\text{22-gram}}(l) = 1$

M2 = |vocab|-gram model = 16-gram model $\rightarrow$ $P_{\text{16-gram}}(l) = \mathcal{L}_{\text{16-gram}}(l) = ?$

$P($ [ The ][ course ][ COMP8730 ][ is ][ about ][ nlp ][ . ][ The ][ instructor ][ 's ][ name ][ is ][ Hossein ][ . ][ There ][ are ][ 13 ][ students ][ in ][ the ][ class ][ . ] $)=$

$P($ [ The ] $)P($ [ course ][ COMP8730 ][ is ][ about ][ nlp ][ . ][ The ][ instructor ][ 's ][ name ][ is ][ Hossein ][ . ][ There ][ are ][ 13 ] .... | [ The ] $)$

$P($ [ The ] $)P($ [ course ] | [ The ] $)P($ [ COMP8730 ][ is ][ about ][ nlp ][ . ][ The ][ instructor ][ 's ][ name ][ is ][ Hossein ][ . ][ There ][ are ] .... | [ The ][ course ] $)$

# Likelihood for a Language Model

$l$: [ The ][ course ][ COMP8730 ][ is ][ about ][ nlp ][ . ][ The ][ instructor ][ 's ][ name ][ is ][ Hossein ][ . ][ There ][ are ][ 13 ][ students ][ in ][ the ][ class ][ . ]

M1 = |token|-gram model = 22-gram model $\rightarrow$ $P_{22\text{-gram}}(l)$ = $\mathcal{L}_{22\text{-gram}}(l)$ = 1

M2 = |vocab|-gram model = 16-gram model $\rightarrow$ $P_{16\text{-gram}}(l)$ = $\mathcal{L}_{16\text{-gram}}(l)$ = ?

$P($ [ The ][ course ][ COMP8730 ][ is ][ about ][ nlp ][ . ][ The ][ instructor ][ 's ][ name ][ is ][ Hossein ][ . ][ There ][ are ][ 13 ][ students ][ in ][ the ][ class ][ . ] $)=$

$P($ [ The ] $)P($ [ course ][ COMP8730 ][ is ][ about ][ nlp ][ . ][ The ][ instructor ][ 's ][ name ][ is ][ Hossein ][ . ][ There ][ are ][ 13 ] .... $|$ [ The ] $)$

$P($ [ The ] $)P($ [ course ] $|$ [ The ] $)P($ [ COMP8730 ][ is ][ about ][ nlp ][ . ][ The ][ instructor ][ 's ][ name ][ is ][ Hossein ][ . ][ There ][ are ] .... $|$ [ The ][ course ] $)$

$P($ [ The ] $)P($ [ course ] $|$ [ The ] $)P($ [ COMP8730 ] $|$ [ The ][ course ] $)P($ [ is ][ about ][ nlp ][ . ] ... There ][ are ] .... $|$ [ The ][ course ][ COMP8730 ] $)$

# Likelihood for a Language Model

$l$: [ The ][ course ][ COMP8730 ][ is ][ about ][ nlp ][ . ][ The ][ instructor ][ 's ][ name ][ is ][ Hossein ][ . ][ There ][ are ][ 13 ][ students ][ in ][ the ][ class ][ . ]

M1 = |token|-gram model = 22-gram model → $P_{22\text{-gram}}(l) = \mathcal{L}_{22\text{-gram}}(l) = 1$

M2 = |vocab|-gram model = 16-gram model → $P_{16\text{-gram}}(l) = \mathcal{L}_{16\text{-gram}}(l) = ?$

P([ The ][ course ][ COMP8730 ][ is ][ about ][ nlp ][ . ][ The ][ instructor ][ 's ][ name ][ is ][ Hossein ][ . ][ There ][ are ][ 13 ][ students ][ in ][ the ][ class ][ . ]) =

P([ The ]) P([ course ][ COMP8730 ][ is ][ about ][ nlp ][ . ][ The ][ instructor ][ 's ][ name ][ is ][ Hossein ][ . ][ There ][ are ][ 13 ] .... | [ The ])

P([ The ]) P([ course ] | [ The ]) P([ COMP8730 ][ is ][ about ][ nlp ][ . ][ The ][ instructor ][ 's ][ name ][ is ][ Hossein ][ . ][ There ][ are ] .... | [ The ][ course ])

P([ The ]) P([ course ] | [ The ]) P([ COMP8730 ] | [ The ][ course ]) P([ is ][ about ][ nlp ][ . ] ... There ][ are ] .... | [ The ][ course ] [ COMP8730 ])

P([ The ]) P([ course ] | [ The ]) P([ COMP8730 ] | [ The ][ course ]) ... P([ name ] | [ The ][ course ] [ COMP8730 ] .... [ 's ]) ... P([ . ] | [The] ... [ class ])

# Likelihood for a Language Model

$l$: [ The ][ course ][ COMP8730 ][ is ][ about ][ nlp ][ . ][ The ][ instructor ][ 's ][ name ][ is ][ Hossein ][ . ][ There ][ are ][ 13 ][ students ][ in ][ the ][ class ][ . ]

M1 = |token|-gram model = 22-gram model → $P_{22\text{-gram}}(l) = \mathcal{L}_{22\text{-gram}}(l) = 1$

M2 = |vocab|-gram model = 16-gram model → $P_{16\text{-gram}}(l) = \mathcal{L}_{16\text{-gram}}(l) = ?$

$P($ [ The ] $)$ ... $P($ [ are ] | [ The ][ course ] [ COMP8730 ] .... [ There ] $)$ ... $P($ [ . ] | [ The ][ course ][ COMP8730 ] .... [ . ] [ The ] ...[ class ] $)$

Only the last 15 words

Only the last 15 words

cannot be considered in 16-gram model

# Likelihood for a Language Model

$l$: [ The ][ course ][ COMP8730 ][ is ][ about ][ nlp ][ . ][ The ][ instructor ][ 's ][ name ][ is ][ Hossein ][ . ][ There ][ are ][ 13 ][ students ][ in ][ the ][ class ][ . ]

M1 = |token|-gram model = 22-gram model $\rightarrow$ P$_{22\text{-gram}}(l) = \mathcal{L}_{22\text{-gram}}(l) = 1$

M2 = |vocab|-gram model = 16-gram model $\rightarrow$ P$_{16\text{-gram}}(l) = \mathcal{L}_{16\text{-gram}}(l) = ?$

M3 = 3-gram model = $\rightarrow$ P$_{3\text{-gram}}(l) = \mathcal{L}_{3\text{-gram}}(l) = ?$

P([ The ])P([ course ] | [ The ])P([ COMP8730 ] | [ The ][course]) ... P([ are ] | [ There ][ is ]) ... P([ . ] | ... [ the ][ class ])

cannot be considered in 3-gram model

# Likelihood for a Language Model

$l$: [ The ][ course ][ COMP8730 ][ is ][ about ][ nlp ][ . ][ The ][ instructor ][ 's ][ name ][ is ][ Hossein ][ . ][ There ][ are ][ 13 ][ students ][ in ][ the ][ class ][ . ]

M1 = |token|-gram model = 22-gram model $\rightarrow$ P$_{22\text{-gram}}(l)$ = $\mathcal{L}_{22\text{-gram}}(l)$ = 1

M2 = |vocab|-gram model = 16-gram model $\rightarrow$ P$_{16\text{-gram}}(l)$ = $\mathcal{L}_{16\text{-gram}}(l)$ = ?

M3 = 3-gram model = $\rightarrow$ P$_{3\text{-gram}}(l)$ = $\mathcal{L}_{3\text{-gram}}(l)$ = ?

M4 = 2-gram model = $\rightarrow$ P$_{2\text{-gram}}(l)$ = $\mathcal{L}_{2\text{-gram}}(l)$ = ?

P($_{[\text{ The }]}$)P($_{[\text{ course }]}$ |$_{[\text{ The }]}$)P($_{[\text{ COMP8730 }]}$|$_{[\text{ The }][\text{course}]}$) … P($_{[\text{ are }]}$ | $_{…\ [\text{ is }]}$) … P($_{[\text{ . }]}$|$_{…\ [\text{ class }]}$)

cannot be considered in 2-gram model

# Likelihood for a Language Model

$l$: [ The ][ course ][ COMP8730 ][ is ][ about ][ nlp ][ . ][ The ][ instructor ][ 's ][ name ][ is ][ Hossein ][ . ][ There ][ are ][ 13 ][ students ][ in ][ the ][ class ][ . ]

M1 = |token|-gram model = 22-gram model → $P_{22\text{-gram}}(l) = \mathcal{L}_{22\text{-gram}}(l) = 1$

M2 = |vocab|-gram model = 16-gram model → $P_{16\text{-gram}}(l) = \mathcal{L}_{16\text{-gram}}(l) = ?$

M3 = 3-gram model = → $P_{3\text{-gram}}(l) = \mathcal{L}_{3\text{-gram}}(l) = ?$

M4 = 2-gram model = → $P_{2\text{-gram}}(l) = \mathcal{L}_{2\text{-gram}}(l) = ?$

M5 = 1-gram model = → $P_{1\text{-gram}}(l) = \mathcal{L}_{1\text{-gram}}(l) = ?$

$P_{([\text{ The }])}P_{([\text{ course }])}P_{([\text{ COMP8730 }])} \dots P_{([\text{ are }])} \dots P_{([\text{ . }])}$ No History!

# Likelihood for a Language Model

$l$: [ The ][ course ][ COMP8730 ][ is ][ about ][ nlp ][ . ][ The ][ instructor ][ 's ][ name ][ is ][
Hossein ][ . ][ There ][ are ][ 13 ][ students ][ in ][ the ][ class ][ . ]

$w\,w\,w\,w\,w\,w = P(w)\,P(w\mid w)\,P(w\mid ww) \cdots \neq P(w)P(w)\cdots P(w)$

M1 = |token|-gram model = 22-gram model $\rightarrow$ $P_{22\text{-gram}}(l) = \mathcal{L}_{22\text{-gram}}(l) = 1$

M2 = |vocab|-gram model = 16-gram model $\rightarrow$ $P_{16\text{-gram}}(l) = \mathcal{L}_{16\text{-gram}}(l) = ?$

M3 = 3-gram model = $\rightarrow$ $P_{3\text{-gram}}(l) = \mathcal{L}_{3\text{-gram}}(l) = ?$
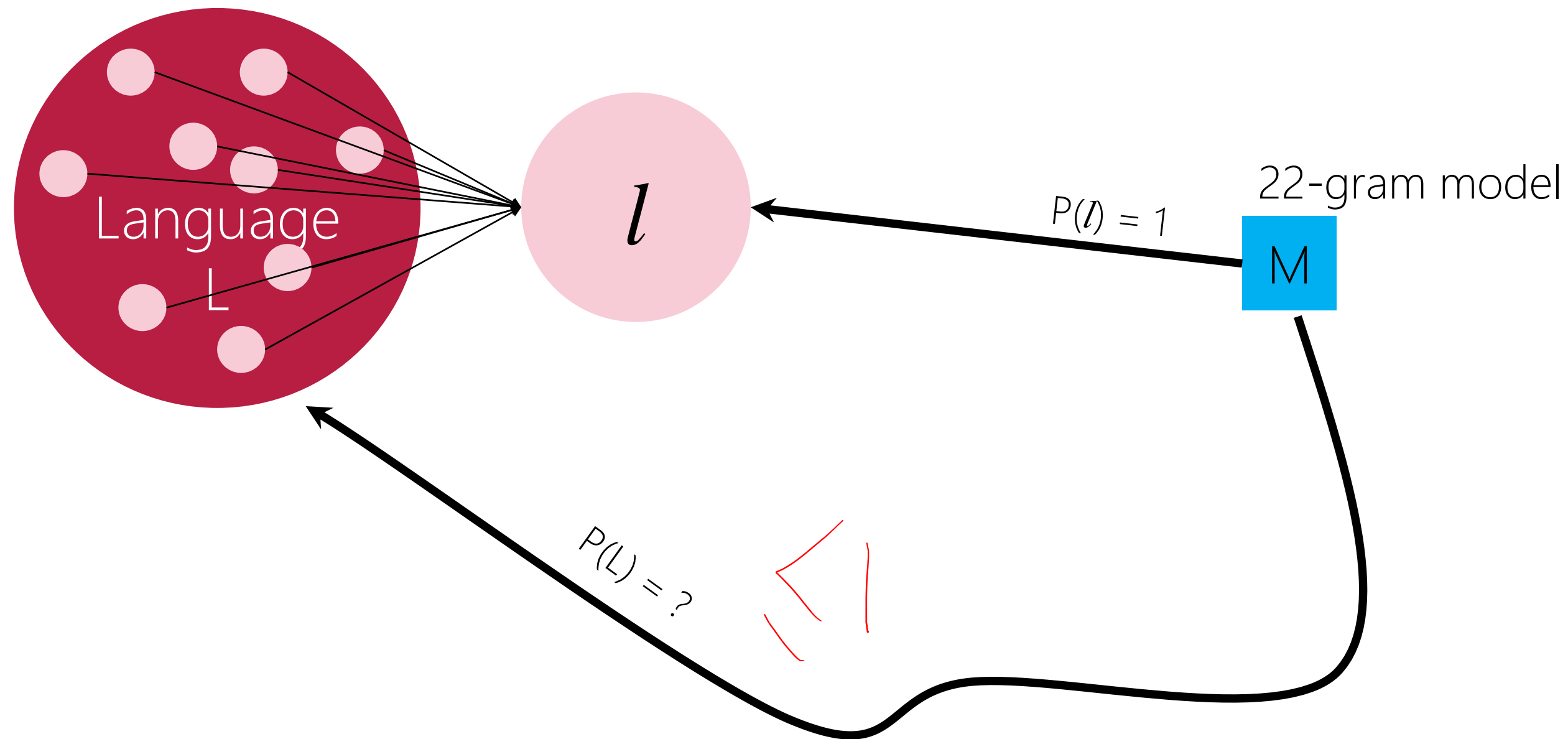
M4 = 2-gram model = $\rightarrow$ $P_{2\text{-gram}}(l) = \mathcal{L}_{2\text{-gram}}(l) = ?$ $\Pi P(w)^n$

M5 = 1-gram model = $\rightarrow$ $P_{1\text{-gram}}(l) = \mathcal{L}_{1\text{-gram}}(l) = ?$

$P(w) = \dfrac{\#\,ww}{\#\,w}$   $P(w\mid w) = \dfrac{\#\,ww}{\#\,w} = \dfrac{4}{5}$  #1

Do you think $\mathcal{L}_{M\text{-}\{2..5\}}(l) \geq \mathcal{L}_{22\text{-gram}}(l)$?

Language L

l

22-gram model

M

$P(l) = 1$

Unseen word, sentence, ...

$P(L) = ?$

# Likelihood for a Language Model

*l*: [ The ][ course ][ COMP8730 ][ is ][ about ][ nlp ][ . ][ The ][ instructor ][ 's ][ name ][ is ][ Hossein ][ . ][ There ][ are ][ 13 ][ students ][ in ][ the ][ class ][ . ]

*l'*: [ Hossein ][ is ][ the ][ name ][ . ]

$P(\text{Hossein}) \times P(\text{is} \mid \text{ho}) \longrightarrow P(\text{the} \mid \text{hossein is})$

M1 = |token|-gram model = 22-gram model $\rightarrow P_{22\text{-gram}}(l') = \mathcal{L}_{22\text{-gram}}(l') = ?$

M2 = |vocab|-gram model = 16-gram model $\rightarrow P_{16\text{-gram}}(l') = \mathcal{L}_{16\text{-gram}}(l') = ?$

M3 = 3-gram model = $\rightarrow P_{3\text{-gram}}(l) = \mathcal{L}_{3\text{-gram}}(l') = ?$

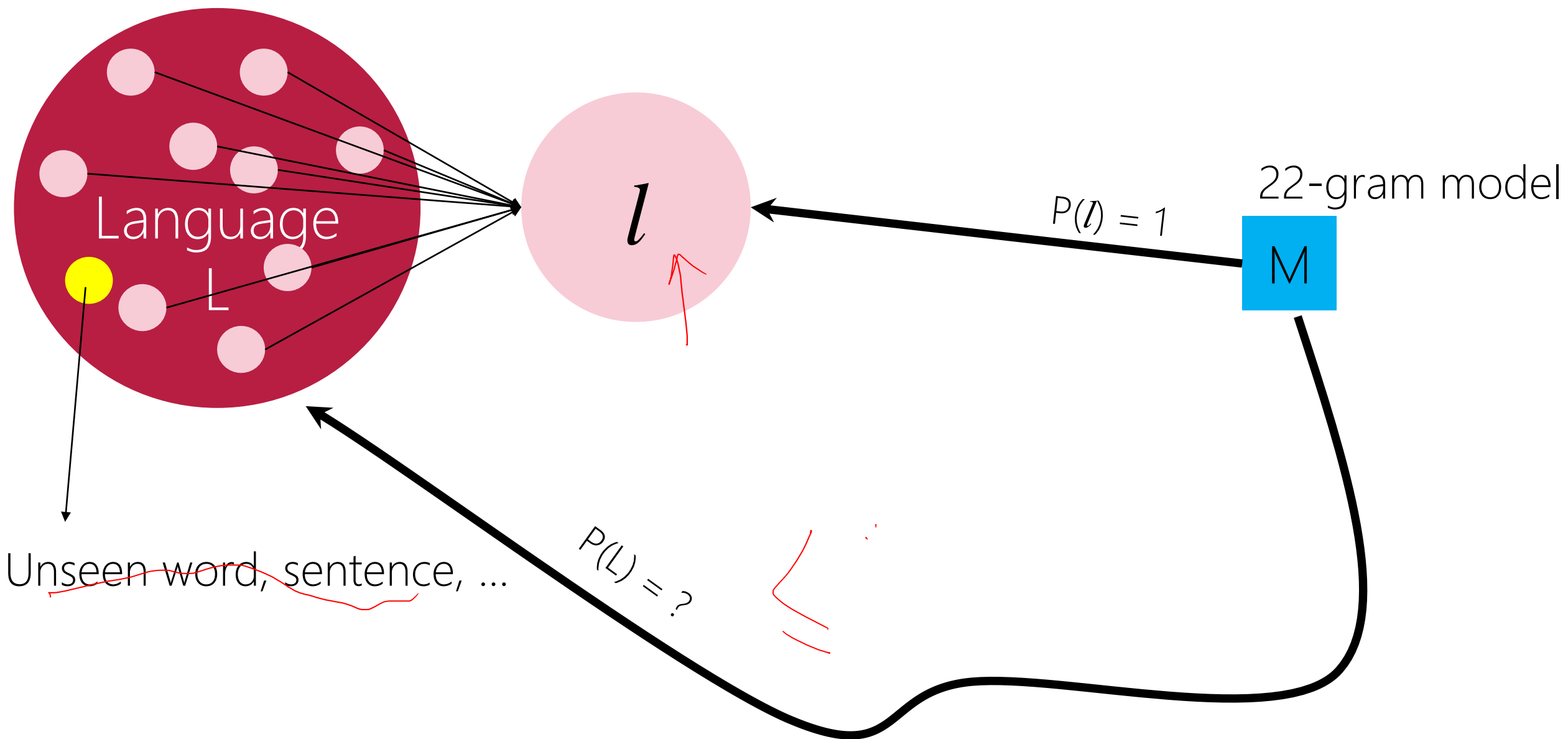M4 = 2-gram model = $\rightarrow P_{2\text{-gram}}(l) = \mathcal{L}_{2\text{-gram}}(l') = ?$

M5 = 1-gram model = $\rightarrow P_{1\text{-gram}}(l) = \mathcal{L}_{1\text{-gram}}(l') = ?$

# Likelihood for a Language Model

$l$: [ The ][ course ][ COMP8730 ][ is ][ about ][ nlp ][ . ][ The ][ instructor ][ 's ][ name ][ is ][ Hossein ][ . ][ There ][ are ][ 13 ][ students ][ in ][ the ][ class ][ . ]

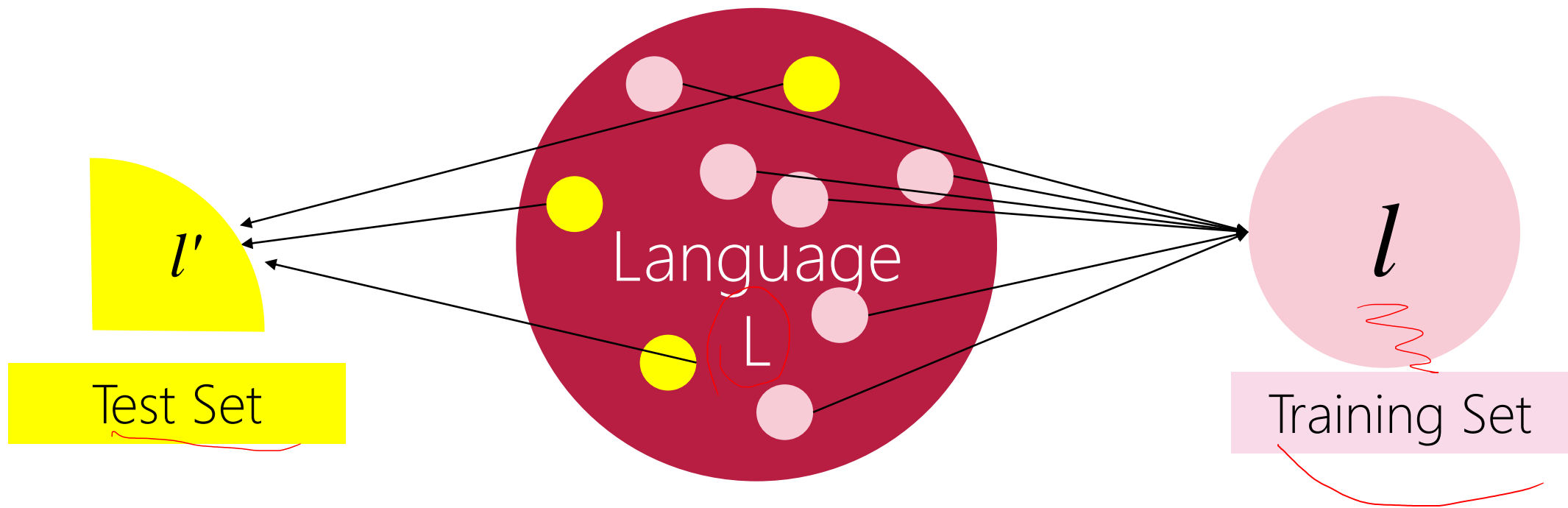$l'$: [ Hossein ][ is ][ the ][ name ][ . ]

M1 = |token|-gram model = 22-gram model $\rightarrow$ $P_{22\text{-gram}}(l') = \mathcal{L}_{22\text{-gram}}(l') = 0$

M2 = |vocab|-gram model = 16-gram model $\rightarrow$ $P_{16\text{-gram}}(l') = \mathcal{L}_{16\text{-gram}}(l') = 0$

M3 = 3-gram model = $\rightarrow$ $P_{3\text{-gram}}(l) = \mathcal{L}_{3\text{-gram}}(l') = 0$

M4 = 2-gram model = $\rightarrow$ $P_{2\text{-gram}}(l) = \mathcal{L}_{2\text{-gram}}(l') = 0$
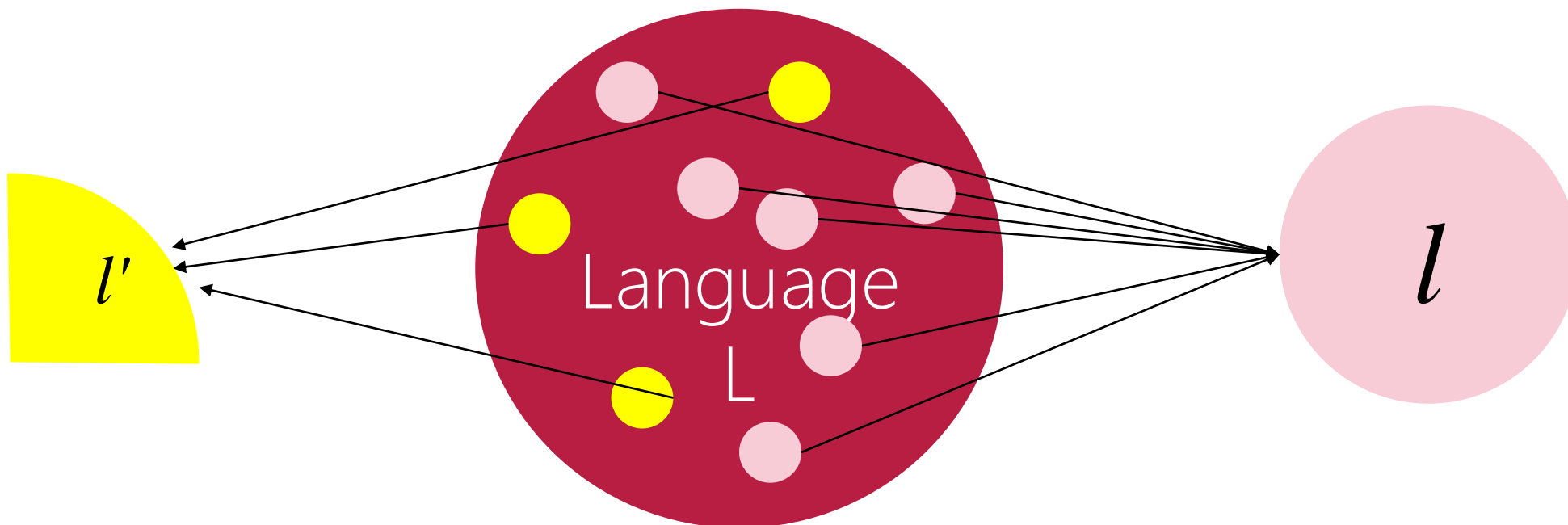
M5 = 1-gram model = $\rightarrow$ $P_{1\text{-gram}}(l) = \mathcal{L}_{1\text{-gram}}(l') = $ Nonzero!

Language L

Test Set $l'$

Training Set $l$

| M is trained on $l$ and does not know anything about $l'$ | | |
|---|---|---|
| | $P_M(l) = \mathcal{L}_M(l)$ | $P_M(l') = \mathcal{L}_M(l')$ |
| M10 | High | High |
| M20 | High | Low |
| M30 | Low | High |
| M40 | Low | Low |

heterarchy

| M is trained on $l$ and does not know anything about $l'$ | | |
|---|---|---|
| | $P_M(l) = \mathcal{L}_M(l)$ | $P_M(l') = \mathcal{L}_M(l')$ |
| M10 | High | High |
| M20 | High | Low |
| M30 | Low | High |
| M40 | Low | Low |

overfitting

?

underfitting

$$\hat{M} = \text{argmax}_{M \in \text{Models}} \mathcal{L}_M(l)$$

1

0

Log Likelihood

$$\hat{M} = \underset{M \in \text{Models}}{\arg\max} \ \log \mathcal{L}_M(l) \quad \substack{0 \\ \uparrow \\ -\infty}$$

# Evaluating Language Models

## Quantitative → Perplexity

How perplexed (confused) a language model is to communicate!
Lower perplexity, the better!

# Perplexity

$$PP_M(l') = \mathcal{L}_M(l')^{\frac{-1}{|l'|}} = \sqrt[|l'|]{\frac{1}{\mathcal{L}(l')}}$$

M is trained on $l$ and test on $l'$

Higher $\mathcal{L}_M(l')$, lower perplexity, the better!

# Perplexity

Unigram approx.: $\sqrt[|l'|]{\dfrac{1}{\mathcal{L}_M(l')}} = \sqrt[|l'|]{\dfrac{1}{\prod_{k=1}^{|l'|} P(w_k)}}$ ; $w_i \in l'$

$$= \sqrt[|l'|]{\dfrac{1}{P(w_k)^{|l'|}}}$$ ; if uniform distribution over words

$$= \dfrac{1}{P(w_k)} = \dfrac{1}{\frac{1}{|V|}} = |V| \text{ the size of vocabs}$$

If LM wants to select a word, it is perplexed in the factor of |V|

# Intuition of Perplexity

- The Shannon Game:
  - How well can we predict the next word?

https://www.youtube.com/watch?v=NCyCkgMLRiY

Lecture 14 — Evaluation and Perplexity — [ NLP || Dan Jurafsky || Stanford University ]

# Perplexity

Unigram approx.: $\sqrt[|l'|]{\dfrac{1}{\mathcal{L}_M(l')}} = \sqrt[|l'|]{\dfrac{1}{\prod_{k=1}^{|l'|} P(w_k)}}$ ; $w_i \in l'$

Bigram approx.: $\sqrt[|l'|]{\dfrac{1}{\mathcal{L}_M(l')}} = \sqrt[|l'|]{\dfrac{1}{\prod_{k=1}^{|l'|} P(w_i|w_{i-1})}}$ ; $w_{i-1}w_i \in l'$

Trigram approx.: $\sqrt[|l'|]{\dfrac{1}{\mathcal{L}_M(l')}} = \sqrt[|l'|]{\dfrac{1}{\prod_{k=1}^{|l'|} P(w_i|w_{i-2}w_{i-1})}}$ ; $w_{i-2}w_{i-1}w_i \in l'$

# Perplexity

$l$ : Wall Street Journal
Size: 38 million words
Vocab (Types): 19,979
$l'$ : 1.5 million words

|  | Unigram | Bigram | Trigram |
|---|---|---|---|
| Perplexity | 962 | 170 | 109 |

Is 4-Gram better?

# Zeros!

# Likelihood for a Language Model

$l$: [ The ][ course ][ COMP8730 ][ is ][ about ][ nlp ][ . ][ The ][ instructor ][ 's ][ name ][ is ][ Hossein ][ . ][ There ][ are ][ 13 ][ students ][ in ][ the ][ class ][ . ]

$l'$: [ Hossein ][ is ][ the ][ name ] [ of ] [ instructor ] [ . ]

M1 = |token|-gram model = 22-gram model $\rightarrow$ $P_{22\text{-gram}}(l') = \mathcal{L}_{22\text{-gram}}(l') = 0$

M2 = |vocab|-gram model = 16-gram model $\rightarrow$ $P_{16\text{-gram}}(l') = \mathcal{L}_{16\text{-gram}}(l') = 0$

M3 = 3-gram model = $\rightarrow$ $P_{3\text{-gram}}(l) = \mathcal{L}_{3\text{-gram}}(l') = 0$

M4 = 2-gram model = $\rightarrow$ $P_{2\text{-gram}}(l) = \mathcal{L}_{2\text{-gram}}(l') = 0$

M5 = 1-gram model = $\rightarrow$ $P_{1\text{-gram}}(l) = \mathcal{L}_{1\text{-gram}}(l') = 0!$ Why?

$P(l') = 0$

# Zeros!

Not all unigrams are available in training set! E.g., [of]

Not all bigrams are available in training set! E.g., [Hossein][is]

Not all trigrams are available in training set! …

# Zeros!

1) Vocabulary + <UNK>  BPE
2) Train Vocabulary + Learn Unseen Tokens (Subwords)
3) Smoothing

# Zeros! <UNK>

## learn the stat of unseen tokens

1) Pick a dictionary $D$
2) From $w \in l$ such that $w \notin D$, (oov) replace it with <UNK>
3) Train model
4) At test, from $w \in l'$, if $w \notin l$ (unseen), replace it with <UNK>

# Zeros! <UNK>
## learn the stat of unseen tokens

*Hossein => I*

*D:* [The][course][is][about][instructor][name][There][are][13][students][in][class]

*l:* ['The][course][<UNK>][is][about][<UNK>][<UNK>][The][instructor][<UNK>][name][is][<UNK>][<UNK>][There][are][13][students][in][the][class][.]

*l':* [<UNK>][is][the][name][<UNK>][<UNK>][<UNK>][<UNK>]

M1 = 22-gram model → P($l'$) = $\mathcal{L}_{22\text{-gram}}(l')$ = 0

M2 = 16-gram model → P($l'$) = $\mathcal{L}_{16\text{-gram}}(l')$ = 0

M3 = 3-gram model = P($l'$) = $\mathcal{L}_{3\text{-gram}}(l')$ = 0

M4 = 2-gram model = P($l'$) = $\mathcal{L}_{2\text{-gram}}(l')$ = 0

M5 = 1-gram model = P($l'$) = $\mathcal{L}_{1\text{-gram}}(l')$ = Nonzero! (Why?)

# Zeros! <UNK>

## learn the stat of unseen tokens

$D$: [The][is]

$l$: ['The][<UNK>][<UNK>][is][<UNK>][<UNK>][<UNK>][The][<UNK>][<UNK>][<UNK>][is][<UNK>][<UNK>][<UNK>][<UNK>]...[<UNK>][the][<UNK>][<UNK>]

$l'$: [<UNK>][is][the][<UNK>][<UNK>][<UNK>][<UNK>][<UNK>]

M1 = 22-gram model ➔ P($l'$) = $\mathcal{L}_{\text{22-gram}}(l')$ = 0

M2 = 16-gram model ➔ P($l'$) = $\mathcal{L}_{\text{16-gram}}(l')$ = 0

M3 = 3-gram model = P($l'$) = $\mathcal{L}_{\text{3-gram}}(l')$ =0

M4 = 2-gram model = P($l'$) = $\mathcal{L}_{\text{2-gram}}(l')$ = Nonzero!

M5 = 1-gram model = P($l'$) = $\mathcal{L}_{\text{1-gram}}(l')$ = Nonzero!
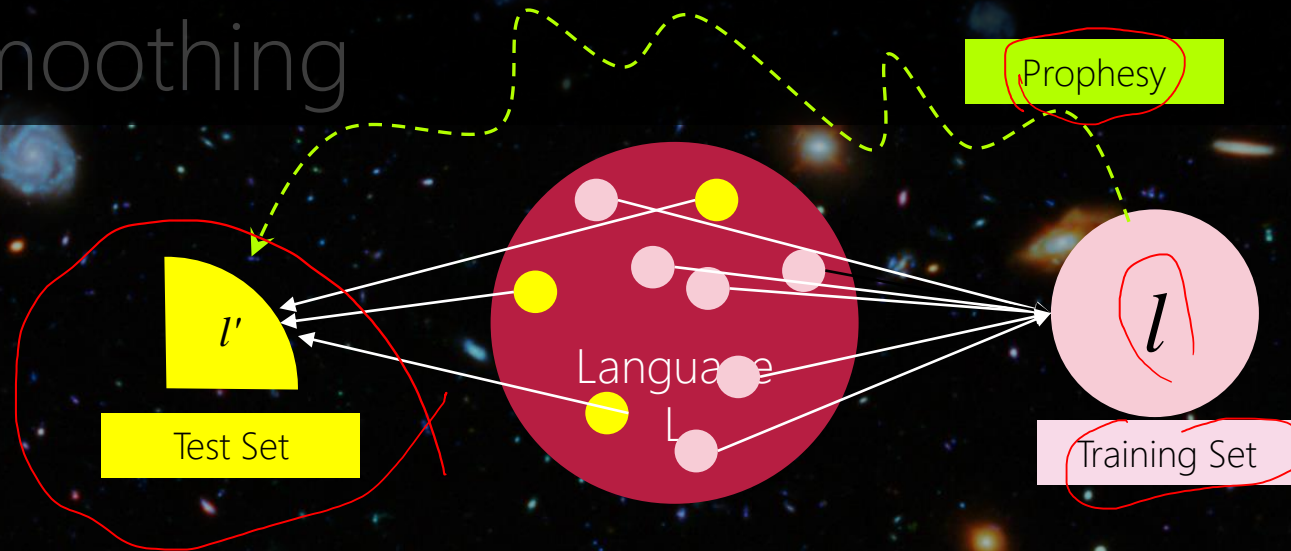
# Zeros! <UNK>

learn the stat of unseen tokens

All the model should generate is stream of <UNK>s!

Pick a Small Dictionary

Gives higher probability (lower perplexity) at test

# Zeros!

1) Vocabulary + <UNK>

2) Train Vocabulary + Learn Unseen Tokens (Subwords)

3) Smoothing

# Zeros!

1) Vocabulary + <UNK>
2) Train Vocabulary + Learn Unseen Tokens (Subwords)
3) Smoothing
   A. Add-1 (Laplace) or Add-k, → k={1,2,...}
   B. Backoff
   C. Interpolation
   D. ...

# Zeros! Add-k

Add *k* unit to all counts, so zero entries become *k*

Add-1 is called Laplace

Unigram model: $P(w_i) = \dfrac{\#w_i + k}{|\text{tokens}| + k \times |\text{vocabs}|}$

Bigram model: $P(w_i|w_{i-1}) = \dfrac{\#(w_{i-1}w_i) + k}{\#(w_{i-1}) + k \times ?}$

Trigram model: $P(w_i|w_{i-2}w_{i-1}) = \dfrac{\#(w_{i-2}w_{i-1}w_i) + k}{\#(w_{i-2}w_{i-1}) + k \times ?}$

D = apple

(the man) 1
(man the) 1

# Zeros! Backoff

if n-Gram have not seen, try (n-1)-Gram

Trigram model: $P(w_i|w_{i-2}w_{i-1}) = \dfrac{\#(w_{i-2}w_{i-1}w_i)}{\#(w_{i-2}w_{i-1})} = 0$

Bigram model: $P(w_i|w_{i-1}) = \dfrac{\#(w_{i-1}w_i)}{\#(w_{i-1})} = 0$

Unigram model: $P(w_i) = \dfrac{\#w_i}{|\text{tokens}|}$

# Zeros! Interpolation

P(n-Gram) is linear interpolation of all (n-i)-Grams: i={1,2,..., n-1}.

Trigram model: $P(w_i|w_{i-2}w_{i-1}) = \lambda_1 \dfrac{\#(w_{i-2}w_{i-1}w_i)}{\#(w_{i-2}w_{i-1})} +$

Bigram model: $P(w_i|w_{i-1}) = \lambda_2 \dfrac{\#(w_{i-1}w_i)}{\#(w_{i-1})} +$

Unigram model: $P(w_i) = \lambda_3 \dfrac{\#w_i}{|\text{tokens}|}$

$\sum \lambda_i = 1$

# Kneser-Ney Smoothing

Kneser, R. and Ney, H. (1995). Improved backing-off for M-gram language modeling. In ICASSP-95, Vol. 1, 181–184.

Chen, S. F. and Goodman, J. (1999). An empirical study of smoothing techniques for language modeling. Computer Speech and Language, 13, 359–394.

# *Evaluating* Language Models

Quantitative → Extrinsic vs. Intrinsic

*Evaluating* Language Models

Quantitative → Spell Correction

$$P(w \mid w_{i+1} \ldots w_{i+n-2}\, w_{i+n-1}) = \frac{P(wi_{+1} \ldots w_{i+n-2}\, w_{i+n-1}\, w)}{P(wi_{+1} \ldots w_{i+n-2}\, w_{i+n-1})}$$

$$= \frac{\#(wi_{+1} \ldots w_{i+n-2}\, w_{i+n-1}\, w)}{\#(wi_{+1} \ldots w_{i+n-2}\, w_{i+n-1})}$$

$$W_{i+1} \ldots W_{i+n-2}\, W_{i+n-1} \longrightarrow W_{i+n}$$

More helpful a language model in finding correct spells, the better!

$$P(w \mid w_{i+1} \ldots w_{i+n-2} \, w_{i+n-1}) = \frac{P(wi_{+1} \ldots w_{i+n-2} \, w_{i+n-1} w)}{P(wi_{+1} \ldots w_{i+n-2} \, w_{i+n-1})}$$

$$= \frac{\#(wi_{+1} \ldots w_{i+n-2} \, w_{i+n-1} w)}{\#(wi_{+1} \ldots w_{i+n-2} \, w_{i+n-1})}$$

$$W_{i+1} \ldots W_{i+n-2} \, W_{i+n-1} \longrightarrow W_{i+n}$$

More helpful a language model in finding correct spells, the better!

Is this judgment correct?