

# GLMTK — Notes

Lukas Schmelzeisen  
lukas@uni-koblenz.de

August 10, 2014

## Contents

<b>1</b>	<b>Attribution</b>	<b>2</b>
<b>2</b>	<b>Notation</b>	<b>2</b>
<b>3</b>	<b><i>n</i>-Gram Probability Estimators</b>	<b>2</b>
3.1	Kinds of Probability Estimators . . . . .	3
3.2	Tests for Probability Estimators . . . . .	3
3.3	Substitute Probability Estimators . . . . .	4
3.4	Fraction Estimators . . . . .	4
3.4.1	MaximumLikelihoodEstimator . . . . .	4
3.4.2	“False”MaximumLikelihoodEstimator . . . . .	4
3.4.3	ContinuationMaximumLikelihoodEstimator . . . . .	5
3.5	Discount Estimators . . . . .	5
3.5.1	AbsoluteDiscountEstimator . . . . .	5
3.6	Backoff Estimators . . . . .	5
3.7	Interpolation Estimators . . . . .	6
3.8	Combination Estimator . . . . .	6
<b>4</b>	<b>TODO</b>	<b>6</b>

# 1 Attribution

I learned about most matter described in this document during my conversations with RENE PICKHARDT. Most of the formulas described here can be attributed to his research.

Besides that, much stems from other language modeling research papers.

## 2 Notation

$\Sigma^n$	Set of all $n$ -Grams
$-$	Skipped Word
$\mathcal{S}_i$	$i$ -th Word in Sequence
$ \mathcal{S} $	Number of Words in Sequence
$c(\mathcal{S})$	Absolute Count of Sequence
$N_{1+}(\mathcal{S})$	Continuation Count of Sequence
$N = c(-)$	Number of Words
$V = N_{1+}(-)$	Vocabulary Size

## 3 $n$ -Gram Probability Estimators

A  $n$ -Gram Probability Estimator is a function  $P : \Sigma^n \rightarrow [0, 1]$  which returns the probability of a  $n$ -Gram *Sequence*  $\mathcal{S}$  for a fixed  $n$ -Gram *History*  $\mathcal{H}$ .

For easier handling we define the *Full Sequence* as the concatenation of history and sequence ( $\mathcal{S}_F = \mathcal{H} * \mathcal{S}$ ), and the *Full History* as the concatenation of history and skipped sequence ( $\mathcal{H}_F = \mathcal{H} * \underbrace{- \dots -}_{|\mathcal{S}| \text{ many}}$ )

An observation on the counts of  $n$ -Grams:

$$c(\mathcal{H}) = 0 \implies c(\mathcal{H}_F) = 0 \implies c(\mathcal{S}_F) = 0 \quad (1)$$

For histories we define the predicate of a (un-)seen history. Note that this defines the empty history as “seen”, which is a choice that was made in order to make the definitions and implementations of estimators more natural.

$$\begin{aligned} \mathcal{H} \text{ seen} &\iff \mathcal{H} = \emptyset \vee c(\mathcal{H}) \neq 0 \\ \mathcal{H} \text{ unseen} &\iff \mathcal{H} \neq \emptyset \wedge c(\mathcal{H}) = 0 \end{aligned} \quad (2)$$

### 3.1 Kinds of Probability Estimators

$n$ -Gram probability estimators can be separated into two categories, according to which mathematical type of probability they implement: *Conditional Probabilities* or *Marginal Probabilities*. How you actually estimate the probability of a sequence depends on what kind of estimator you are using.

$n$ -Gram probability with conditional probability estimators:

$$P(w_1^n) = P(w_n | w_1^{n-1}) \cdot P(w_{n-1} | w_1^{n-2}) \cdots P(w_1) \quad (3)$$

$n$ -Gram probability with marginal probability estimators:

$$P(w_1^n) = P(w_n | w_1^{n-1}) \cdot P(w_{n-1} | w_1^{n-2}) \cdots P(w_1) \quad (4)$$

Conditional and marginal probabilities differ on how they handle the case of an unseen history. Conditional probabilities have  $P(\mathcal{S} | \mathcal{H} \text{ unseen}) = 0$  while marginal probabilities have  $P(\mathcal{S} | \mathcal{H} \text{ unseen}) = P_{\text{Substitute}}(\mathcal{S} | \mathcal{H})$ .

### 3.2 Tests for Probability Estimators

In order for probability estimators to be probability measures, the following equations / tests should hold:

NGramProbabilitiesSumTest:

$$\sum_{\mathcal{S} \in \Sigma^n} P(\mathcal{S}) = 1 \quad (5)$$

FixedHistoryProbabilitiesSumTest:

Conditional:

$$\begin{aligned} \forall \mathcal{H} \in \Sigma^n : (\mathcal{H} \text{ seen} &\implies \sum_{\mathcal{S} \in \Sigma} P(\mathcal{S} | \mathcal{H}) = 1) \wedge \\ (\mathcal{H} \text{ unseen} &\implies \sum_{\mathcal{S} \in \Sigma} P(\mathcal{S} | \mathcal{H}) = 0) \end{aligned} \quad (6)$$

Marginal:

$$\forall \mathcal{H} \in \Sigma^n : \sum_{\mathcal{S} \in \Sigma} P(\mathcal{S} | \mathcal{H}) = 1 \quad (7)$$

### 3.3 Substitute Probability Estimators

Substitute Probability Estimators are used in a context where other probability estimators cannot use their usual algorithm to estimate the probability of a sequence. They then instead use  $P_{\text{Substitute}}$  to calculate that probability.

Let  $P_{\text{Substitute}} \in \{P_{\text{Uniform}}, P_{\text{AbsUnigram}}, P_{\text{ContUnigram}}\}$  fixed globally at program start.

$$P_{\text{Uniform}}(\mathcal{S}|\mathcal{H}) = \frac{1}{V} \quad (8)$$

$$P_{\text{AbsUnigram}}(\mathcal{S}|\mathcal{H}) = \frac{c(\mathcal{S}_1)}{N} \quad (9)$$

$$P_{\text{ContUnigram}}(\mathcal{S}|\mathcal{H}) = \frac{N_{1+}(-\mathcal{S}_1)}{N_{1+}(-)} \quad (10)$$

All substitute probability estimators are marginal probability estimators.

### 3.4 Fraction Estimators

Fraction Estimators are probability estimators that have the form  $\frac{n}{d}$ .

Conditional:

$$P_{\text{Frac}[n,d]}(\mathcal{S}|\mathcal{H}) = \begin{cases} 0 & \text{if } \mathcal{H} \text{ unseen} \vee d = 0 \\ \frac{n}{d} & \text{else} \end{cases} \quad (11)$$

Marginal:

$$P_{\text{Frac}[n,d]}(\mathcal{S}|\mathcal{H}) = \begin{cases} P_{\text{Substitute}}(\mathcal{S}|\mathcal{H}) & \text{if } \mathcal{H} \text{ unseen} \vee d = 0 \\ \frac{n}{d} & \text{else} \end{cases} \quad (12)$$

#### 3.4.1 MaximumLikelihoodEstimator

$$P_{\text{MLE}}(\mathcal{S}|\mathcal{H}) = P_{\text{Frac}[c(\mathcal{S}_F), c(\mathcal{H}_F)]}(\mathcal{S}|\mathcal{H}) \quad (13)$$

#### 3.4.2 “False”MaximumLikelihoodEstimator

FMLE only works in the marginal probability setting.

$$P_{\text{FMLE}}(\mathcal{S}|\mathcal{H}) = P_{\text{Frac}[c(\mathcal{S}_F), c(\mathcal{H})]}(\mathcal{S}|\mathcal{H}) \quad (14)$$

### 3.4.3 ContinuationMaximumLikelihoodEstimator

$$P_{\text{CMLE}}(\mathcal{S}|\mathcal{H}) = P_{\text{Frac}[N_{1+}(-\mathcal{S}_F), N_{1+}(-\mathcal{H}_F)]}(\mathcal{S}|\mathcal{H}) \quad (15)$$

## 3.5 Discount Estimators

A Discount Estimator takes any kind of fraction estimator and subtracts some discount from the numerator, in order to free probability mass to be used for smoothing. Obviously this means, that discount estimators are no longer probability estimators, and will not pass tests. Instead they have to be used in conjunction with an Interpolation Estimator. Discount estimators are still fraction estimators though.

$$P_{\text{Discount}[D, P_{\text{Frac}}[n, d]]}(\mathcal{S}|\mathcal{H}) = P_{\text{Frac}[\max(0, n-D), d]}(\mathcal{S}|\mathcal{H}) \quad (16)$$

With  $D : \mathcal{H} \rightarrow [0, 1]$ .

### 3.5.1 AbsoluteDiscountEstimator

$$P_{\text{AbsDiscount}[D, P_{\text{Frac}}[n, d]]}(\mathcal{S}|\mathcal{H}) = P_{\text{Discount}[D, P_{\text{Frac}}[n, d]]}(\mathcal{S}|\mathcal{H}) \quad (17)$$

With  $D \in [0, 1]$ .

## 3.6 Backoff Estimators

Conditional:

$$P_{\text{Backoff}[P_\alpha, P_\beta]}(\mathcal{S}|\mathcal{H}) = \begin{cases} 0 & \text{if } \mathcal{H} = \emptyset \\ \gamma(\mathcal{H}) \cdot P_\beta(\mathcal{S}|\hat{\mathcal{H}}) & \text{if } c(\mathcal{S}_F) = 0 \\ P_\alpha(\mathcal{S}|\mathcal{H}) & \text{else} \end{cases} \quad (18)$$

Marginal:

$$P_{\text{Backoff}[P_\alpha, P_\beta]}(\mathcal{S}|\mathcal{H}) = \begin{cases} P_{\text{Substitute}}(\mathcal{S}|\mathcal{H}) & \text{if } \mathcal{H} = \emptyset \\ \gamma(\mathcal{H}) \cdot P_\beta(\mathcal{S}|\hat{\mathcal{H}}) & \text{if } c(\mathcal{S}_F) = 0 \\ P_\alpha(\mathcal{S}|\mathcal{H}) & \text{else} \end{cases} \quad (19)$$

With  $P_\alpha, P_\beta$  any probability estimators and *backoff coefficient*  $\gamma$ :

$$\gamma(\mathcal{H}) = \frac{1 - \sum_{\mathcal{S} \in \Sigma: c(\mathcal{H}\mathcal{S}) > 0} P_\alpha(\mathcal{S}|\mathcal{H})}{\sum_{\mathcal{S} \in \Sigma: c(\mathcal{H}\mathcal{S}) = 0} P_\beta(\mathcal{S}|\hat{\mathcal{H}})} \quad (20)$$

### 3.7 Interpolation Estimators

$$P_{\text{Interpol}}[P_{\text{Discount}}[D, P_{\text{Frac}}[n, d]], P_{\beta}](\mathcal{S}|\mathcal{H}) = \begin{cases} P_{\text{Frac}}[n, d](\mathcal{S}|\mathcal{H}) & \text{if } \mathcal{H} = \emptyset \vee \mathcal{H} \text{ contains only } - \\ P_{\text{Discount}}[D, P_{\text{Frac}}[n, d]](\mathcal{S}|\mathcal{H}) + \gamma(D, d, H) \cdot P_{\beta}(\mathcal{S}|\mathcal{H}) & \text{else} \end{cases} \quad (21)$$

With the  $P_{\beta}$  any probability estimator and *interpolation coefficient*  $\gamma$ :

$$\gamma(D, d, H) = \begin{cases} 0 & \text{if } d = 0 \\ \frac{D \cdot N_{1+}(\mathcal{H} -)}{d} & \text{else} \end{cases} \quad (22)$$

### 3.8 Combination Estimator

A Combination Estimator mixes two other probability estimators.

$$P_{\text{Comb}}[\lambda, P_{\alpha}, P_{\beta}](\mathcal{S}|\mathcal{H}) = \lambda \cdot P_{\alpha}(\mathcal{S}|\mathcal{H}) + (1 - \lambda) \cdot P_{\beta}(\mathcal{S}|\mathcal{H}) \quad (23)$$

With  $P_{\alpha}, P_{\beta}$  any probability estimators and  $\lambda \in [0, 1]$ .

## 4 TODO

- Interpol estimator still fails test for conditional case.
- FMLE still doesn't pass any tests.
- How to make `FixedHistoryProbabilitiesSumTest` with Continuation Estimators not be a mess?
- How to do `FixedHistoryProbabilitiesSumTest` with Combination Estimator.