

Supplementary Material for Language Modeling with Power Low Rank Ensembles

Ankur P. Parikh, Avneesh Saluja, Chris Dyer, Eric P. Xing
Carnegie Mellon University

[]

The primary purpose of the supplementary material is to provide a proof of Lemma 4. We also show that Lemma 4 extends to $n > 2$.

1 Proof of Lemma 4

Lemma 1. Let $P_{\text{plre}}(w_i|w_{i-1})$ indicate the PLRE smoothed conditional probability and $\widehat{P}(w)$ indicate the MLE probability of w . Then,

$$\widehat{P}(w) = \sum_{w_{i-1}} P_{\text{plre}}(w_i|w_{i-1}) \widehat{P}(w_{i-1}) \quad (1)$$

Proof. Assume the following more general form where multiple low rank matrices can be used i.e.:

$$P_{\text{plre}}(w_i|w_{i-1}) = P_{D_0}^{\text{alt}}(w_i|w_{i-1}) + \gamma_0(w_{i-1}) \left(Z_{D_1}^{(\rho_1, \kappa_1)}(w_i|w_{i-1}) + \dots \right. \\ \left. + \gamma_{\eta-1}(w_{i-1}) \left(Z_{D_\eta}^{(\rho_\eta, \kappa_\eta)}(w_i|w_{i-1}) + \gamma_\eta(w_{i-1}) \left(Z^{(\rho_{\eta+1}=0, \kappa_{\eta+1}=1)}(w_i|w_{i-1}) \right) \right) \dots \right) \quad (2)$$

where we note that $Z^{(\rho_{\eta+1}=0, \kappa_{\eta+1}=1)}(w_i|w_{i-1})$ is equivalent to $P^{\text{alt}}(w_i)$. It is assumed that $1 \geq \rho_0 \geq \dots \geq \rho_{\eta+1} = 0$.

First unroll the recursion and rewrite $P_{\text{plre}}(w_i|w_{i-1})$ as:

$$P_{\text{plre}}(w_i|w_{i-1}) = \sum_{j=0}^{\eta+1} \gamma_{0:j-1}(w_{i-1}) Z_{D_j}^{(\rho_j, \kappa_j)}(w_i|w_{i-1})$$

where $\gamma_{0:j-1}(w_{i-1}) = \prod_{h=0}^{j-1} \gamma_h(w_{i-1})$ and $\gamma_{0:-1}(w_{i-1}) = 1$. Note that $P_{\text{pwr}}(w_i|w_{i-1})$ can be written in the same way.

$$P_{\text{pwr}}(w_i|w_{i-1}) = \sum_{j=0}^{\eta} \gamma_{0:j}(w_{i-1}) Y_{D_j}^{(\rho_j)}(w_i|w_{i-1}) \quad (3)$$

Note that $P_{\text{pwr}}(w_i|w_{i-1})$ already satisfies the marginal constraint i.e.

$$\widehat{P}(w) = \sum_{w_{i-1}} P_{\text{pwr}}(w_i|w_{i-1}) \widehat{P}(w_{i-1}) \quad (4)$$

because the discounts were chosen such that $P_{\text{pwr}}(w_i|w_{i-1}) = \widehat{P}(w_i|w_{i-1})$

Thus it suffices to show that for all $j = 0, \dots, \eta+1$:

$$\sum_{w_{i-1}} \widehat{P}(w_{i-1}) \gamma_{0:j-1}(w_{i-1}) Y_{D_j}^{(\rho_j)}(w_i|w_{i-1}) = \sum_{w_{i-1}} \widehat{P}(w_{i-1}) \gamma_{0:j-1}(w_{i-1}) Z_{D_j}^{(\rho_j, \kappa_j)}(w_i|w_{i-1}) \quad (5)$$

The statement above is trivially true when $j = 0$. For all other cases, note that due to the way we have set the discounts, $\gamma_{0:j-1}$ takes a special form:

$$\begin{aligned} \prod_{h=0}^{j-1} \gamma_h(w_{i-1}) &= \frac{d_* \sum_i c_{i,i-1}^{\rho_1}}{\sum_i c_{i,i-1}^{\rho_0}} \frac{d_* \sum_i c_{i,i-1}^{\rho_2}}{\sum_i c_{i,i-1}^{\rho_1}} \dots \frac{d_* \sum_i c_{i,i-1}^{\rho_j}}{\sum_i c_{i,i-1}^{\rho_{j-1}}} \\ &= \frac{(d_*)^j \sum_i c_{i,i-1}^{\rho_j}}{\sum_i c_{i,i-1}} \end{aligned} \quad (6)$$

Using this form in Eq. 5 and simplifying yields:

$$\sum_{w_{i-1}} \left(\sum_i c_{i,i-1}^{\rho_j} \right) \mathbf{Y}_{\mathbf{D}_j}^{(\rho_j)}(w_i | w_{i-1}) = \sum_{w_{i-1}} \left(\sum_i c_{i,i-1}^{\rho_j} \right) \mathbf{Z}_{\mathbf{D}_j}^{(\rho_j, \kappa_j)}(w_i | w_{i-1})$$

which is equivalent to requiring that

$$\sum_{w_{i-1}} \mathbf{Y}_{\mathbf{D}_j}^{(\rho_j)}(w_i, w_{i-1}) = \sum_{w_{i-1}} \mathbf{Z}_{\mathbf{D}_j}^{(\rho_j, \kappa_j)}(w_i, w_{i-1}) \quad (7)$$

which holds because rank minimization under gKL preserves row and column sums. \square

2 Generalization to $n > 2$

Theorem 1. Let $P_{\text{plre}}(w_i | w_{i-n+1}^{i-1})$ indicate the PLRE smoothed conditional probability and $\widehat{P}(w)$ indicate the MLE probability of w . Then,

$$\widehat{P}(w) = \sum_{w_{i-n+1}^{i-1}} P_{\text{plre}}(w_i | w_{i-n+1}^{i-1}) \widehat{P}(w_{i-n+1}^{i-1}) \quad (8)$$

Proof. Recall that,

$$\begin{aligned} P_{\text{plre}}(w_i | w_{i-n+1}^{i-1}) &= P_{\mathbf{D}_0}^{\text{alt}}(w_i | w_{i-n+1}^{i-1}) \\ &+ \gamma_0(w_{i-n+1}^{i-1}) \left(\mathbf{Z}_{\mathbf{D}_1}^{(\rho_1, \kappa_1)}(w_i | w_{i-n+1}^{i-1}) + \dots \right. \\ &+ \gamma_{\eta-1}(w_{i-n+1}^{i-1}) \left(\mathbf{Z}_{\mathbf{D}_\eta}^{(\rho_\eta, \kappa_\eta)}(w_i | w_{i-n+1}^{i-1}) \right. \\ &\left. \left. + \gamma_\eta(w_{i-n+1}^{i-1}) \left(P_{\text{plre}}(w_i | w_{i-n+2}^{i-1}) \right) \right) \dots \right) \end{aligned} \quad (9)$$

Define,

$$\begin{aligned} P_{\text{pwr}}(w_i | w_{i-n+1}^{i-1}) &= P_{\mathbf{D}_0}^{\text{alt}}(w_i | w_{i-n+1}^{i-1}) \\ &+ \gamma_0(w_{i-n+1}^{i-1}) \left(\mathbf{Y}_{\mathbf{D}_1}^{(\rho_1, \kappa_1)}(w_i | w_{i-n+1}^{i-1}) + \dots \right. \\ &+ \gamma_{\eta-1}(w_{i-n+1}^{i-1}) \left(\mathbf{Y}_{\mathbf{D}_\eta}^{(\rho_\eta, \kappa_\eta)}(w_i | w_{i-n+1}^{i-1}) \right. \\ &\left. \left. + \gamma_\eta(w_{i-n+1}^{i-1}) \left(P_{\text{pwr}}(w_i | w_{i-n+2}^{i-1}) \right) \right) \dots \right) \end{aligned} \quad (10)$$

where with a little abuse of notation

$$\mathbf{Y}_{\mathbf{D}_j}^{\rho_j}(w_i | w_{i-n'+1}^{i-1}) = \frac{\tilde{c}(w_i, w_{i-n'+1}^{i-1})^{\rho_j} - \mathbf{D}_j(w_i, w_{i-n'+1}^{i-1})}{\sum_{w_i} \tilde{c}(w_i, w_{i-n'+1}^{i-1})^{\rho_j}} \quad (11)$$

and

$$\tilde{c}(w_i, w_{i-n'+1}^{i-1}) = \begin{cases} c(w_i, w_{i-n'+1}^{i-1}), & \text{if } n' = n \\ N_-(w_{i-n'+1}^i) & \text{if } n' < n \end{cases}$$

Furthermore, define

$$\begin{aligned} P_{\text{pwr}}^{\text{terms}}(w_i | w_{i-n'+1}^{i-1}) &= P_{D_0}^{\text{alt}}(w_i | w_{i-n'+1}^{i-1}) \\ &+ \gamma_0(w_{i-n'+1}^{i-1}) \left(\mathbf{Y}_{D_1}^{(\rho_1, \kappa_1)}(w_i | w_{i-n'+1}^{i-1}) + \dots \right. \\ &\quad \left. + \gamma_{n-1}(w_{i-n'+1}^{i-1}) \left(\mathbf{Y}^{(\rho_n, \kappa_n)}(w_i | w_{i-n'+1}^{i-1}) \right) \dots \right) \end{aligned} \quad (12)$$

Note that because of the way the discounts are computed in Algorithm 1,

$$P_{\text{pwr}}^{\text{terms}}(w_i | w_{i-n'+1}^{i-1}) = P^{\text{alt}}(w_i | w_{i-n'+1}^{i-1}) \quad (13)$$

for all $n' \leq n$.

As a result, (for some choice of discount)

$$P_{\text{pwr}}(w_i | w_{i-n+1}^{i-1}) = P_{\text{kn}}(w_i | w_{i-n+1}^{i-1}) \quad (14)$$

Since, we know that Kneser Ney satisfies the marginal constraint (Chen and Goodman, 1999) this implies that,

$$\widehat{P}(w) = \sum_{w_{i-n+1}^{i-1}} P_{\text{pwr}}(w_i | w_{i-n+1}^{i-1}) \widehat{P}(w_{i-n+1}^{i-1}) \quad (15)$$

Thus, all we have to do is prove that

$$\sum_{w_{i-n+1}^{i-1}} P_{\text{pwr}}(w_i | w_{i-n+1}^{i-1}) \widehat{P}(w_{i-n+1}^{i-1}) = \sum_{w_{i-n+1}^{i-1}} P_{\text{pre}}(w_i | w_{i-n+1}^{i-1}) \widehat{P}(w_{i-n+1}^{i-1}) \quad (16)$$

Now, we follow the same argument as with $n = 2$ (i.e. unrolling the recursion and applying the fact that gKL preserves row/column sums).

For notational simplicity assume that $n = 3$. Then, we can write $P_{\text{pwr}}(w_i | w_{i-n+1}^{i-1})$ as:

$$P_{\text{pwr}}(w_i | w_{i-2}^{i-1}) = \sum_{j=0}^{\eta} \gamma_{0:j-1}(w_{i-2}^{i-1}) \mathbf{Y}_{D_j}^{(\rho_j, \kappa_j)}(w_i | w_{i-2}^{i-1}) + \sum_{j=0}^{\eta+1} \gamma_{0:\eta}(w_{i-2}^{i-1}) \gamma_{0:j-1}(w_{i-1}) \mathbf{Y}_{D_j}^{(\rho_j, \kappa_j)}(w_i | w_{i-1}) \quad (17)$$

where $\gamma_{0:-1}(w_{i-2}^{i-1}) = 1$ and

$$\begin{aligned} \gamma_{0:j-1}(w_{i-2}^{i-1}) &= \prod_{h=0}^{j-1} \gamma_h(w_{i-2}^{i-1}) = \frac{d_* \sum_i \tilde{c}_{i,i-1,i-2}^{\rho_1}}{\sum_i \tilde{c}_{i,i-1,i-2}^{\rho_0}} \frac{d_* \sum_i \tilde{c}_{i,i-1,i-2}^{\rho_2}}{\sum_i \tilde{c}_{i,i-1,i-2}^{\rho_1}} \dots \frac{d_* \sum_i \tilde{c}_{i,i-1,i-2}^{\rho_j}}{\sum_i \tilde{c}_{i,i-1,i-2}^{\rho_{j-1}}} \\ &= \frac{(d_*)^j \sum_i \tilde{c}_{i,i-1,i-2}^{\rho_j}}{\sum_i \tilde{c}_{i,i-1,i-2}} \end{aligned} \quad (18)$$

Here $\tilde{c}_{i,i-1,i-2}$ is shorthand for $\tilde{c}(w_i, w_{i-2}^{i-1})$.

Similarly, $\gamma_{0:-1}(w_{i-1}) = 1$ and

$$\begin{aligned}\gamma_{0:j-1}(w_{i-1}) &= \prod_{h=0}^{j-1} \gamma_h(w_{i-1}) = \frac{d_* \sum_i \tilde{c}_{i,i-1}^{\rho_1}}{\sum_i \tilde{c}_{i,i-1}^{\rho_0}} \frac{d_* \sum_i \tilde{c}_{i,i-1}^{\rho_2}}{\sum_i \tilde{c}_{i,i-1}^{\rho_1}} \dots \frac{d_* \sum_i \tilde{c}_{i,i-1}^{\rho_j}}{\sum_i \tilde{c}_{i,i-1}^{\rho_{j-1}}} \\ &= \frac{(d_*)^j \sum_i \tilde{c}_{i,i-1}^{\rho_j}}{\sum_i \tilde{c}_{i,i-1}}\end{aligned}\quad (19)$$

(Again, it is assumed that $1 \geq \rho_0 \geq \dots \geq \rho_{\eta+1} = 0$.)

Analogously,

$$P_{\text{pre}}(w_i | w_{i-2}^{i-1}) = \sum_{j=0}^{\eta} \gamma_{0:j-1}(w_{i-2}^{i-1}) \mathbf{Z}_{D_j}^{(\rho_j, \kappa_j)}(w_i | w_{i-2}^{i-1}) + \sum_{j=0}^{\eta+1} \gamma_{0:\eta}(w_{i-2}^{i-1}) \gamma_{0:j-1}(w_{i-1}) \mathbf{Z}_{D_j}^{(\rho_h, \kappa_j)}(w_i | w_{i-1}) \quad (20)$$

Now for any trigram term we prove that

$$\sum_{w_{i-2}^{i-1}} \gamma_{0:j-1}(w_{i-2}^{i-1}) \mathbf{Y}_{D_j}^{(\rho_j, \kappa_j)}(w_i | w_{i-2}^{i-1}) \widehat{P}(w_{i-2}^{i-1}) = \sum_{w_{i-2}^{i-1}} \gamma_{0:j-1}(w_{i-2}^{i-1}) \mathbf{Z}_{D_j}^{(\rho_j, \kappa_j)}(w_i | w_{i-2}^{i-1}) \widehat{P}(w_{i-2}^{i-1}) \quad (21)$$

Plugging in the definition of $\gamma_{0:j-1}$ and simplifying gives

$$\sum_{w_{i-2}^{i-1}} \left(\sum_i \tilde{c}_{i,i-1,i-2}^{\rho_j} \right) \mathbf{Y}_{D_j}^{(\rho_j, \kappa_j)}(w_i | w_{i-2}^{i-1}) = \sum_{w_{i-2}^{i-1}} \left(\sum_i \tilde{c}_{i,i-1,i-2}^{\rho_j} \right) \mathbf{Z}_{D_j}^{(\rho_j, \kappa_j)}(w_i | w_{i-2}^{i-1}) \quad (22)$$

which is equivalent to

$$\sum_{w_{i-2}^{i-1}} \mathbf{Y}_{D_j}^{(\rho_j, \kappa_j)}(w_i, w_{i-2}^{i-1}) = \sum_{w_{i-2}^{i-1}} \mathbf{Z}_{D_j}^{(\rho_j, \kappa_j)}(w_i, w_{i-2}^{i-1}) \quad (23)$$

which holds because of the definition of \mathbf{Z} and the fact that rank minimization under gKL preserves row/column sums.

Now consider any bigram term. We seek to show that:

$$\begin{aligned}\sum_{w_{i-2}^{i-1}} \gamma_{0:\eta}(w_{i-2}^{i-1}) \gamma_{0:j-1}(w_{i-1}) \mathbf{Y}_{D_j}^{(\rho_j, \kappa_j)}(w_i | w_{i-1}) \widehat{P}(w_{i-2}^{i-1}) \\ = \sum_{w_{i-2}^{i-1}} \gamma_{0:\eta}(w_{i-2}^{i-1}) \gamma_{0:j-1}(w_{i-1}) \mathbf{Z}_{D_j}^{(\rho_j, \kappa_j)}(w_i | w_{i-1}) \widehat{P}(w_{i-2}^{i-1})\end{aligned}\quad (24)$$

Substituting definition of $\gamma_{0:\eta}(w_{i-2}^{i-1})$ gives

$$\begin{aligned}\sum_{w_{i-2}^{i-1}} \frac{(d_*)^{\eta+1} \sum_i \tilde{c}_{i,i-1,i-2}^{\rho_{\eta+1}}}{\sum_i \tilde{c}_{i,i-1,i-2}} \gamma_{0:j-1}(w_{i-1}) \mathbf{Y}_{D_j}^{(\rho_j, \kappa_j)}(w_i | w_{i-1}) \widehat{P}(w_{i-2}^{i-1}) \\ = \sum_{w_{i-2}^{i-1}} \frac{(d_*)^{\eta+1} \sum_i \tilde{c}_{i,i-1,i-2}^{\rho_{\eta+1}}}{\sum_i \tilde{c}_{i,i-1,i-2}} \gamma_{0:j-1}(w_{i-1}) \mathbf{Z}_{D_j}^{(\rho_j, \kappa_j)}(w_i | w_{i-1}) \widehat{P}(w_{i-2}^{i-1})\end{aligned}\quad (25)$$

Simplifying and pushing in the sum over w_{i-2} gives,

$$\sum_{w_{i-1}} \left(\sum_{i,i-2} \tilde{c}_{i,i-1,i-2}^{\rho_{\eta+1}} \right) \gamma_{0:j-1}(w_{i-1}) \mathbf{Y}_{D_j}^{(\rho_j, \kappa_j)}(w_i | w_{i-1}) = \sum_{w_{i-1}} \left(\sum_{i,i-2} \tilde{c}_{i,i-1,i-2}^{\rho_{\eta+1}} \right) \gamma_{0:j-1}(w_{i-1}) \mathbf{Z}_{D_j}^{(\rho_j, \kappa_j)}(w_i | w_{i-1}) \quad (26)$$

Note that since $\rho_{\eta+1} = 0$, $\sum_{i,i-2} \tilde{c}_{i,i-1,i-2}^{\rho_{\eta+1}=0} = \sum_i \tilde{c}_{i,i-1}$ (by definition of \tilde{c}).

Using this fact and substituting definition of $\gamma_{0:j-1}(w_{i-1})$ gives

$$\sum_{w_{i-1}} \left(\sum_i \tilde{c}_{i,i-1} \right) \frac{(d_*)^j \sum_i \tilde{c}_{i,i-1}^{\rho_j}}{\sum_i \tilde{c}_{i,i-1}} \mathbf{Y}_{D_j}^{(\rho_j, \kappa_j)}(w_i | w_{i-1}) = \sum_{w_{i-1}} \left(\sum_i \tilde{c}_{i,i-1} \right) \frac{(d_*)^j \sum_i \tilde{c}_{i,i-1}^{\rho_j}}{\sum_i \tilde{c}_{i,i-1}} \mathbf{Z}_{D_j}^{(\rho_j, \kappa_j)}(w_i | w_{i-1}) \quad (27)$$

Simplifying gives,

$$\sum_{w_{i-1}} \mathbf{Y}_{D_j}^{(\rho_j)}(w_i, w_{i-1}) = \sum_{w_{i-1}} \mathbf{Z}_{D_j}^{(\rho_j, \kappa_j)}(w_i, w_{i-1}) \quad (28)$$

which holds because rank minimization under KL divergence preserves row and column sums. \square

References

- Stanley F. Chen and Joshua Goodman. 1999. An empirical study of smoothing techniques for language modeling. *Computer Speech & Language*, 13(4):359–393.