Lecture Series on Predictive Language Models

El Mehdi ISSOUANI

Post-doc: Laboratory of Applied Mathematics of Compiègne (LMAC)
University of Technology of Compiègne (UTC)

Thesis obtained in June 2023 in MODAL'X at University Paris Nanterre

Thursday 6th February, 2025

Lecture 8: Mathematical Formalization of POS Tagging: Maximum Entropy and Extension to Empirical Likelihood

Outline

- NLP Tasks
 - Tokenization
 - Part-of-Speech tagging
- Mathematical Formalization of POS Tagging
 - Notations and Vocabulary
 - Feature-based models
- Maximum entropy principle (MaxEnt)
 - Link to divergences
 - MaxEnt solution
 - MaxEnt for POS Tagging
- Penalized generalized empirical likelihood (PGEL)
 - ullet ϕ^* -discrepancy and Dual Theorem
 - Penalizing the dual program
 - The explicit expression of parameters for POS Tagging

- NLP Tasks
 - Tokenization
 - Part-of-Speech tagging
- Mathematical Formalization of POS Tagging
- Maximum entropy principle (MaxEnt)
- Penalized generalized empirical likelihood (PGEL)

Natural Language Processing

- Tokenization
- Part Of Speech Tagging (POS tagging)

- NLP Tasks
 - Tokenization
 - Part-of-Speech tagging
- Mathematical Formalization of POS Tagging
- Maximum entropy principle (MaxEnt)
- Penalized generalized empirical likelihood (PGEL)

Natural Language Processing

Tokenization examples

[He called Mr. Green at 2 p.m. in St. Louis, Mr. White did not answer. He then left him a voice mail message.]

Sentence tokenization

[He called Mr. Green at 2 p.m. in St. Louis, Mr. White did not answer.] [He then left him a voice mail message.]

Word tokenization

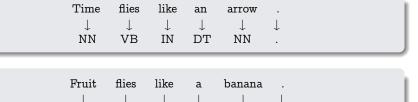
[He], [then], [left], [him], [a], [voice], [mail], [message], [.]

- NLP Tasks
 - Tokenization
 - Part-of-Speech tagging
- 2 Mathematical Formalization of POS Tagging
- Maximum entropy principle (MaxEnt)
- Penalized generalized empirical likelihood (PGEL)

5 / 23

Natural Language Processing

POS tagging examples



NN

VB: Verb - DT: Determiner - NN: Noun singular - IN: Preposition - NNS: Noun plural - JJ: Adjective

JJ NNS VB DT

- NLP Tasks
- 2 Mathematical Formalization of POS Tagging
 - Notations and Vocabulary
 - Feature-based models
- 3 Maximum entropy principle (MaxEnt)
- Penalized generalized empirical likelihood (PGEL)

Notations and Vocabulary

ullet We denote a sentence s with N_s words

$$s = w_1, \ldots, w_i, \ldots, w_{N_s}$$

where w_i represents the *i*th word and $N_s \in \mathbb{N}^*$ is random.

- t_i is the tag of the *i*th word with $t_i \in T = \{NN, NNS, VB, ...\}$ representing the tagset of finite size (example: 7, 36, 87)
- n: Number of words in the whole dataset $n = \sum_{s \in S} N_s$

	w_1	W_2		w_{N-1}	w_N	
	1	1	\uparrow	\uparrow	↑	\uparrow
$\overline{\mathit{sent}} o$	Time	flies	like	an	arrow	•
	1	1	1	1	1	1
$\textit{tags} \rightarrow$	NN	VB	IN	DT	NN	•
	\downarrow	\downarrow	\downarrow	\downarrow	\downarrow	\downarrow
	t_1	t_2		t_{N-1}	t_N	

Mathematical Models for POS Tagging

• Goal : we aim to obtain the best tag sequence

$$(t_{1}^{*},...,t_{N}^{*}) = \underset{(t_{1},...,t_{N}) \in \mathbf{T^{N}}}{\arg \max} \left[p\left(t_{1},...,t_{N}|w_{1},...,w_{N}\right) \right]$$

• Denote $x_1, ..., x_N$ the contexts of each word.

Example:
$$x_i = [w_{i-2}, w_{i-1}, w_i, w_{i+1}, w_{i+2}]$$

We now can write

$$p(t_1,...,t_N|w_1,...,w_N) = \prod_{i=1}^{N} p(t_i|x_i)$$

Mathematical Models for POS Tagging

Goal: we aim to obtain the best tag sequence

$$(t_{1}^{*},...,t_{N}^{*}) = \underset{(t_{1},...,t_{N}) \in \mathbf{T}^{\mathbf{N}}}{\arg\max} \left[p\left(t_{1},...,t_{N}|w_{1},...,w_{N}\right) \right]$$

• Denote $x_1, ..., x_N$ the contexts of each word.

Example:
$$x_i = [w_{i-2}, w_{i-1}, w_i, w_{i+1}, w_{i+2}]$$

We now can write

$$p(t_1,...,t_N|w_1,...,w_N) = \prod_{i=1}^{N} p(t_i|x_i)$$

- NLP Tasks
- 2 Mathematical Formalization of POS Tagging
 - Notations and Vocabulary
 - Feature-based models
- 3 Maximum entropy principle (MaxEnt)
- 4 Penalized generalized empirical likelihood (PGEL)

Feature-based models

 Definition: A feature encodes the binary link between the tag and the context (or word environment).



		 (flies, NNS)	(flies, VB)	• • •	(Time, VB)	• • •
$f(x_2, NNS)$	0	 1	0		0	0
$f(x_2, VB)$	0	 0	1		1	0
$f(x_2, NN)$	0	 0	0		0	0
:	:	:	:	:	:	:

- NLP Tasks
- 2 Mathematical Formalization of POS Tagging
- Maximum entropy principle (MaxEnt)
 - Link to divergences
 - MaxEnt solution
 - MaxEnt for POS Tagging
- Penalized generalized empirical likelihood (PGEL)

Maximum entropy principle (MaxEnt)

A method to infer a measure p(z) defined on a given set $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$ under some constraint \mathcal{P} (Csiszar (1996)[5], Hanson (2012)[6]).

$$\rho^* = \argmax_{\rho \in \mathcal{P}} \{\mathcal{H}(\rho)\} = \argmax_{\rho \in \mathcal{P}} \left\{ -\int \rho(z) \log \rho(z) l(dz) \right\},$$

where
$$\mathcal{P} = \left\{ p : \int f_k(z) p(z) l(dz) = \mu_k, k = 1, \dots, q \right\}.$$

- NLP Tasks
- 2 Mathematical Formalization of POS Tagging
- Maximum entropy principle (MaxEnt)
 - Link to divergences
 - MaxEnt solution
 - MaxEnt for POS Tagging
- 4 Penalized generalized empirical likelihood (PGEL)

Link to divergences

If we have access to a given default distribution $p_0 \in \mathcal{P}$

$$\begin{split} p^* &= \operatorname*{arg\;max} \left\{ \mathcal{H}(p) \right\} = & \operatorname*{arg\;max} \left\{ \mathcal{H}(p_0) - D\left(p, p_0\right) \right\} \\ &= & \operatorname*{arg\;min} \left\{ D\left(p, p_0\right) \right\} \\ &\underset{p \in \mathcal{P}}{=} \end{split}$$

where
$$D(p, p_0) = \int \left[p(z) \log \frac{p(z)}{p_0(z)} - p(z) + p_0(z) \right] l(dz).$$

- NLP Tasks
- 2 Mathematical Formalization of POS Tagging
- Maximum entropy principle (MaxEnt)
 - Link to divergences
 - MaxEnt solution
 - MaxEnt for POS Tagging
- 4 Penalized generalized empirical likelihood (PGEL)

MaxEnt solution

Consider λ_i 's the Kuhn & Tucker coefficients in the optimization program with each λ_i corresponding to a constraint μ_i ,

$$\rho^{*}(z) = \frac{\exp\left(\sum_{k=1,..,q} \lambda_{k} f_{k}(z)\right)}{\int \exp\left(\sum_{k=1,..,q} \lambda_{k} f_{k}(u)\right) l(du)}.$$

MaxEnt solution

Consider λ_i 's the Kuhn & Tucker coefficients in the optimization program with each λ_i corresponding to a constraint μ_i ,

$$\rho^{*}(z) = \frac{\exp\left(\sum_{k=1,..,q} \lambda_{k} f_{k}(z)\right)}{\int \exp\left(\sum_{k=1,..,q} \lambda_{k} f_{k}(u)\right) l(du)}.$$

 \longrightarrow Notice that p(z) doesn't depend on constraints μ_k .

- NLP Tasks
- 2 Mathematical Formalization of POS Tagging
- Maximum entropy principle (MaxEnt)
 - Link to divergences
 - MaxEnt solution
 - MaxEnt for POS Tagging
- Penalized generalized empirical likelihood (PGEL)

Back to POS Tagging

By putting z = (x, t) we can rewrite

Log-linear models, the join probability is given by

$$p^{*}(z) = p^{*}(x,t) = p^{*}(t|x)p^{*}(x) = \frac{\exp\left(\sum_{k=1,..,q} \lambda_{k} f_{k}(x,t)\right)}{\int \exp\left(\sum_{k=1,..,q} \lambda_{k} f_{k}(u,v)\right) l(du,dv)},$$

Conditional probabilities

$$\rho^*(t|x) = \frac{\exp\left(\sum_{k=1,..,q} \lambda_k f_k\left(x,t\right)\right)}{\int \exp\left(\sum_{k=1,..,q} \lambda_k f_k\left(x,t'\right)\right) l'(dt')}.$$

- NLP Tasks
- 2 Mathematical Formalization of POS Tagging
- Maximum entropy principle (MaxEnt)
- Penalized generalized empirical likelihood (PGEL)
 - φ*-discrepancy and Dual Theorem
 - Penalizing the dual program
 - The explicit expression of parameters for POS Tagging

ullet Consider $(\mathcal{Z},\mathcal{A},\mathcal{M})$ where \mathcal{M} is a space of signed measures. For every convex function ϕ , its Fenchel-Legendre transform is given by

$$\varphi^*(z) = \sup_{z \in \mathbb{R}} \{zy - \varphi(z)\}, \ \forall \ y \in \mathbb{R}.$$

13 / 23

ullet Consider $(\mathcal{Z},\mathcal{A},\mathcal{M})$ where \mathcal{M} is a space of signed measures. For every convex function ϕ , its Fenchel-Legendre transform is given by

$$\phi^*(z) = \sup_{z \in \mathbb{R}} \{zy - \phi(z)\}, \ \forall \ y \in \mathbb{R}.$$

ullet Consider the following assumptions for the function ϕ .

13 / 23

• Consider $(\mathcal{Z}, \mathcal{A}, \mathcal{M})$ where \mathcal{M} is a space of signed measures. For every convex function ϕ , its Fenchel-Legendre transform is given by

$$\varphi^*(z) = \sup_{z \in \mathbb{R}} \{zy - \varphi(z)\}, \ \forall \ y \in \mathbb{R}.$$

• Consider the following assumptions for the function φ .

H1 φ is strictly convex and $d(\varphi)$ contains a neighborhood of 0;

ullet Consider $(\mathcal{Z},\mathcal{A},\mathcal{M})$ where \mathcal{M} is a space of signed measures. For every convex function ϕ , its Fenchel-Legendre transform is given by

$$\phi^*(z) = \sup_{z \in \mathbb{R}} \{zy - \phi(z)\}, \ \forall \ y \in \mathbb{R}.$$

• Consider the following assumptions for the function φ .

H1 φ is strictly convex and $d(\varphi)$ contains a neighborhood of 0;

 $\mbox{H2}\ \phi$ is twice differentiable on a neighborhood of 0 ;

ullet Consider $(\mathcal{Z},\mathcal{A},\mathcal{M})$ where \mathcal{M} is a space of signed measures. For every convex function ϕ , its Fenchel-Legendre transform is given by

$$\phi^*(z) = \sup_{z \in \mathbb{R}} \{zy - \phi(z)\}, \ \forall \ y \in \mathbb{R}.$$

- Consider the following assumptions for the function φ .
- **H1** φ is strictly convex and $d(\varphi)$ contains a neighborhood of 0;
- **H2** φ is twice differentiable on a neighborhood of 0;
- **H3** (renormalization) $\phi(0)=0$ and $\phi^{(1)}(0)=0$, $\phi^{(2)}(0)>0$, which implies that ϕ has an unique minimum at zero ;

ullet Consider $(\mathcal{Z},\mathcal{A},\mathcal{M})$ where \mathcal{M} is a space of signed measures. For every convex function ϕ , its Fenchel-Legendre transform is given by

$$\phi^*(z) = \sup_{z \in \mathbb{R}} \{zy - \phi(z)\}, \ \forall \ y \in \mathbb{R}.$$

- Consider the following assumptions for the function φ .
- **H1** φ is strictly convex and $d(\varphi)$ contains a neighborhood of 0;
- **H2** φ is twice differentiable on a neighborhood of 0;
- **H3** (renormalization) $\phi(0)=0$ and $\phi^{(1)}(0)=0$, $\phi^{(2)}(0)>0$, which implies that ϕ has an unique minimum at zero ;
- **H4** φ is differentiable on $d(\varphi)$, that is to say differentiable on $int\{d(\varphi)\}$, with right and left limits on the respective endpoints of the support of $d(\varphi)$, where $int\{.\}$ is the topological interior.

- NLP Tasks
- 2 Mathematical Formalization of POS Tagging
- Maximum entropy principle (MaxEnt)
- Penalized generalized empirical likelihood (PGEL)
 - ullet ϕ^* -discrepancy and Dual Theorem
 - Penalizing the dual program
 - The explicit expression of parameters for POS Tagging

ϕ^* -discrepancy and Dual Theorem

$$I_{\varphi^*}(\mathbb{Q},\mathbb{P}) = \left\{ egin{array}{ll} \int_{\mathcal{Z}} \phi^* \left(rac{d\mathbb{Q}}{d\mathbb{P}} - 1
ight) d\mathbb{P} & ext{if } \mathbb{Q} \ll \mathbb{P} \\ + \infty & ext{else}. \end{array}
ight.$$

Theorem (Borwein & Lewis (1992)[2], Keziou (2003)[7], Bertail (2007)[1])

Let ϕ be a function satisfying assumptions **H1-H3**. If the following qualification constraint holds,

$$\mathsf{Qual}(\mathbb{P}): \begin{cases} \exists \mathbb{T} \in \mathcal{M}, \mathbb{T}f = \mu \text{ and} \\ \inf d(\phi^*) \ < \ \inf_{\mathcal{Z}} \frac{d\mathbb{T}}{d\mathbb{P}} \ \leq \ \sup_{\mathcal{Z}} \frac{d\mathbb{T}}{d\mathbb{P}} \ < \ \sup d(\phi^*) \ \mathbb{P}-\textit{a.s.}, \end{cases}$$

then, we have the dual equality:

$$\inf_{\mathbb{Q}\in\mathcal{M}}\left\{I_{\varphi^*}(\mathbb{Q},\mathbb{P})|\ (\mathbb{Q}-\mathbb{P})f=\mu\right\}=\sup_{\lambda\in\mathbb{R}^q}\left\{\lambda'\mu-\int_{\mathcal{Z}}\varphi(\lambda'f)d\mathbb{P}\right\}$$

ϕ^* -discrepancy and Dual Theorem

$$I_{\varphi^*}(\mathbb{Q},\mathbb{P}) = \left\{ egin{array}{ll} \int_{\mathcal{Z}} \varphi^* \left(rac{d\mathbb{Q}}{d\mathbb{P}} - 1
ight) d\mathbb{P} & ext{if } \mathbb{Q} \ll \mathbb{P} \\ +\infty & ext{else}. \end{array}
ight.$$

Theorem (Borwein & Lewis (1992)[2], Keziou (2003)[7], Bertail (2007)[1])

Let ϕ be a function satisfying assumptions **H1-H3**. If the following qualification constraint holds,

$$\mathsf{Qual}(\mathbb{P}): \begin{cases} \exists \mathbb{T} \in \mathcal{M}, \mathbb{T}f = \mu \text{ and} \\ \inf d(\phi^*) \ < \ \inf_{\mathcal{Z}} \frac{d\mathbb{T}}{d\mathbb{P}} \ \leq \ \sup_{\mathcal{Z}} \frac{d\mathbb{T}}{d\mathbb{P}} \ < \ \sup d(\phi^*) \ \mathbb{P}-\textit{a.s.}, \end{cases}$$

then, we have the dual equality:

$$\beta_{n}(\mu) = \inf_{\mathbb{Q} \in \mathcal{M}} \{I_{\phi^{*}}(\mathbb{Q}, \mathbb{P}) | (\mathbb{Q} - \mathbb{P})f = \mu\} = \sup_{\lambda \in \mathbb{R}^{q}} \left\{ \lambda' \mu - \int_{\mathcal{Z}} \phi(\lambda' f) d\mathbb{P} \right\}$$

GEL solution

If ϕ satisfies H4, then the infimum on the left hand side at \mathbb{Q}^* is given by

$$p^* = \mathbb{Q}^* = (1 + \varphi^{(1)}(\lambda^{*\prime}f))\mathbb{P}.$$

Define the empirical measure: $\mathbb{P}_n = \frac{1}{n} \sum_{i=1}^n \delta_{Z_i}$.

Example of Kullback divergence for measures:

It's the case where $\varphi_0(x) = -x - \log(1-x)$ and $\varphi_0^*(x) = x - \log(1+x)$.

Kullback divergence

$$I_{\varphi_{\mathbf{0}}^*}(\mathbb{Q},\mathbb{P}) = K(\mathbb{Q},\mathbb{P}) = -\int_{\mathcal{Z}} \log(\frac{d\mathbb{Q}}{d\mathbb{P}}) d\mathbb{P} + \int_{\mathcal{Z}} (d\mathbb{Q} - d\mathbb{P}),$$

$$p_i^* = \frac{1}{n(1 - \lambda^{*'}(Z_i - \mu))}$$
 and $\beta_n(\mu) = -1 - \sum_{i=1}^{n} \log(np_i^*) + \sum_{i=1}^{n} p_i^*$.

where,
$$p_i^* = \frac{1}{n(1 - \lambda^{*'}(Z_i - \mu))}$$
.

Examples of relative entropy and χ^2 divergences

• Relative entropy: The particular case of $\varphi_1(x) = e^x - 1 - x$ whose convex conjugate is given by $\varphi_1^*(x) = (x+1)\log(1+x) - x$. Then,

Relative entropy

$$I_{\phi_{\mathbf{1}}^*}(\mathbb{Q},\mathbb{P}) = \left\{ \begin{array}{ll} \int_{\mathcal{Z}} \frac{d\mathbb{Q}}{d\mathbb{P}} \log(\frac{d\mathbb{Q}}{d\mathbb{P}}) d\mathbb{P} - \int_{\mathcal{Z}} (d\mathbb{Q} - d\mathbb{P}) & \text{ if } \mathbb{Q} \ll \mathbb{P} \\ +\infty & \text{ else.} \end{array} \right.$$

And that the optimal weights are given by

$$p_i^* = rac{1}{n} \exp\left(\lambda^{*\prime} f(Z_i, \mu)
ight), \quad ext{where } \lambda^* \mathop{\sim}_{n o \infty} - S_n^{-2} \overline{f}_n$$

• Program solution with χ^2 divergence:

$$\beta_{n}\left(\mu\right)=n\overline{f}_{n}^{\prime}S_{n}^{-2}\overline{f}_{n}\quad\text{where }\overline{f}_{n}=\frac{1}{n}\sum_{i=1}^{n}f(Z_{i},\mu).$$

- NLP Tasks
- 2 Mathematical Formalization of POS Tagging
- Maximum entropy principle (MaxEnt)
- Penalized generalized empirical likelihood (PGEL)
 - φ*-discrepancy and Dual Theorem
 - Penalizing the dual program
 - The explicit expression of parameters for POS Tagging

Penalizing the dual program

- Remember that for the POS Tagging: $n \ll q$
- Penalizing the dual (Chang et al. (2018) [4], Shi (2016) [8])

$$P_{n}(\mu,\lambda) = \mathbb{P}_{n}\left(-\lambda'\left(f\left(x,t\right) - \mu\right) - \varphi(\lambda'\left(f\left(x,t\right) - \mu\right))\right) - \frac{1}{2}||\lambda||_{R}^{2},$$

In the case of Relative entropy divergence:

$$\Longrightarrow P_n(\mu,\lambda) = 1 + \frac{1}{n} \sum_{i=1}^n \left(-\exp(\lambda' \left(f\left(x_i,t_i\right) - \mu\right) \right) - \frac{1}{2} ||\lambda||_R^2 \right).$$

When $R = \rho_n I_q$, then we have asymptotically for λ close to 0,

$$P_n(\mu,\lambda) \approx \frac{1}{n} \sum_{i=1}^n \left(-\lambda' \left(f\left(x_i, t_i\right) - \mu\right) \right) - \frac{1}{2} \lambda' (S_n^2 + \rho_n I_q) \lambda \right)$$

- NLP Tasks
- 2 Mathematical Formalization of POS Tagging
- Maximum entropy principle (MaxEnt)
- Penalized generalized empirical likelihood (PGEL)
 - \bullet ϕ^* -discrepancy and Dual Theorem
 - Penalizing the dual program
 - The explicit expression of parameters for POS Tagging

The explicit expression of λ^* for POS Tagging

The maximum is attained at

$$\begin{split} \lambda_n^* \underset{\infty}{\sim} &- (S_n^2 + \rho_n I_q)^{-1} \mathbb{P}_n(f - \mu) \quad (\text{recall that } \mathbb{P}_n = \frac{1}{n} \sum_i \delta_{Z_i}) \\ \underset{\infty}{\sim} &- (S_n^2 + \rho_n I_q)^{-1} \frac{1}{n} \sum_{i=1}^n \left(f\left(x_i, t_i\right) - \mu \right) \end{split}$$

In the penalized case we can see that the optimal weights depend on $\boldsymbol{\mu}$

Without using the maximum of likelihood (Issouani (2023))

$$\hat{p}\left(t_{i}|x\right) = \frac{e^{-(\tilde{f}_{n}-\mu)'\left(S_{n}^{2}+\rho_{n}I_{q}\right)^{-2}\left(f\left(x,t_{i}\right)-\mu\right)}}{\sum\limits_{t_{k}\in\mathcal{T}} e^{-(\tilde{f}_{n}-\mu)'\left(S_{n}^{2}+\rho_{n}I_{q}\right)^{-2}\left(f\left(x,t_{k}\right)-\mu\right)}},$$

where we recall that $S_n^2 = \frac{1}{n} \sum_{i=1}^n (f(x_i, t_i) - \mu) (f(x_i, t_i) - \mu)'$.

Application for POS Tagging

- PennTreebank corpus: 3914 sentences having a total of 100676 tokens (12408 tokens without repetitions) and 36 tags.
 - \bullet We estimate μ using the empirical mean estimated using the entire initial dataset
 - Then, we split the dataset into a training set (75% of the initial dataset) and a test set

Results:

- Estimation accuracy of 98% (over the training sample)
- Prediction accuracy of 95% (on average over the test samples).

19 / 23

Perspectives

- The impact of divergence choice on probabilities.
- ullet Choose another norm than L_2 for the dual penalizing.



Jonathan M Borwein and Adrian S Lewis. Partially finite convex programming, part ii: Explicit lattice models. *Mathematical Programming*, 57:49–83, 1992.

Michel Broniatowski and Amor Keziou.

Minimization of φ-divergences on sets of signed measures.

Studia Scientiarum Mathematicarum Hungarica, 43(4):403–442, 2006.

Jinyuan Chang, Cheng Yong Tang, and Tong Tong Wu.

A new scope of penalized empirical likelihood with high-dimensional estimating equations.

The Annals of Statistics, 46(6B):3185-3216, 2018.

Imre Csiszár.

Maxent, mathematics, and information theory.

In Maximum entropy and Bayesian methods, pages 35-50. Springer, 1996.



Kenneth M Hanson and Richard N Silver.

Maximum Entropy and Bayesian Methods: Santa Fe, New Mexico, USA, 1995 Proceedings of the Fifteenth International Workshop on Maximum Entropy and Bayesian Methods, volume 79.

Springer Science & Business Media, 2012.



Amor Keziou.

Dual representation of φ -divergences and applications. Comptes Rendus Mathematique, 336(10):857–862, 2003.



Zhentao Shi.

Econometric estimation with high-dimensional moment equalities.

Journal of Econometrics, 195(1):104–119, 2016.

Thank you for your attention!