

Тематическое моделирование образовательных целей пользователей в системе дистанционного образования

Потапова Полина

Московский физико-технический институт
Факультет управления и прикладной математики
Кафедра интеллектуальных систем

Научный руководитель д.ф.-м.н. К. В. Воронцов

Москва,
2020 г.

Актуальность

Пользователи систем дистанционного образования получают персонализированные рекомендации. Для построения ИОТ необходимо учитывать их образовательные цели. Нужно уметь оценивать конкретность целей.

SMART методология постановки целей

- Specific: Конкретная
- Measurable: Измеримая
- Achievable or Attainable: Достижимая
- Relevant: Значимая
- Time bound: Ограниченная по времени

Цель работы

Построить модель бинарной классификации для определения конкретности образовательных целей, написанных на естественном языке.

Хорошо сформулированы относительно SMART

- Изучить историю дизайна, колористику, в совершенстве владеть графическими программами линейки Adobe: Ps, Ai, Ae, Xd, Lr, Id, Pr, а также историю и инструменты менеджмента для поступления в магистратуру Политехнического университета Петра Великого до июня 2020 года.
- Войти в категорию A1 по классификации призывников, что даст возможность попасть в Кремлёвский полк. Deadline: 2021 г.

Плохо сформулированы относительно SMART

- Счастья всем даром, и пусть никто не уйдёт обиженным!
- Хочу быть способным поддержать диалог на любую тему.

Постановка задачи оценки конкретности образовательной цели

Дано

Выборка $X = (x_i, y_i)_{i=1}^l$,

x_i - признаковое описание образовательной цели (документа),
 $y_i \in \{0, 1\}$ - метка класса.

Найти

Алгоритм бинарной классификации $a(x) : X \rightarrow Y = \{0, 1\}$, предсказывающий, является ли образовательная цель конкретной.

Постановка задачи оценки конкретности образовательной цели

Варианты алгоритма $a(x)$

baseline:	Наивный байесовский классификатор, логистическая регрессия, простая тематическая модель классификации
основное решение:	Тематическая модель конкретности образовательных целей

Критерий

Площадь под ROC-кривой (roc_auc), accuracy, precision и recall.

Вероятностная модель порождения текстовой коллекции документов:

$$p(w|d) = \sum_{t \in T} p(w|t)p(t|d) = \sum_{t \in T} \phi_{wt}\theta_{td}, \quad d \in D, w \in W,$$

где $p(w|t) = \phi_{wt}$ – распределение термина $w \in W$ в теме $t \in T$,

$p(t|d) = \theta_{td}$ – распределение темы $t \in T$ в документе $d \in D$.

Принцип максимума регуляризованного правдоподобия:

$$\sum_{w \in W} \sum_{d \in D} n_{dw} \ln \sum_{t \in T} \phi_{wt}\theta_{td} + R(\Phi, \Theta) \rightarrow \max_{\Phi, \Theta}, \quad \sum_{w \in W} \phi_{wt} = 1, \phi_{wt} \geq 0, \sum_{t \in T} \theta_{td} = 1, \theta_{td} \geq 0$$

EM-алгоритм:

$$p_{tdw} = \mathit{norm}_{t \in T}(\phi_{wt}\theta_{td});$$

$$\phi_{wt} = \mathit{norm}_{w \in W} \left(n_{wt} + \phi_{wt} \frac{\partial R}{\partial \phi_{wt}} \right);$$

$$n_{wt} = \sum_{d \in D} n_{dw} p_{tdw};$$

$$\theta_{td} = \mathit{norm}_{t \in T} \left(n_{td} + \theta_{td} \frac{\partial R}{\partial \theta_{td}} \right);$$

$$n_{td} = \sum_{w \in W} n_{dw} p_{tdw};$$

Документ $d \in D$ состоит из токенов w_m разных модальностей $m \in M$.
Принцип максимума регуляризованного правдоподобия:

$$\sum_{m \in M} \tau_m \sum_{d \in D} \sum_{w \in W_m} n_{dw} \ln \sum_{t \in T} \phi_{wt} \theta_{td} + \sum_i \tau_i R_i(\Phi, \Theta) \rightarrow \max_{\Phi, \Theta}$$

Условия неотрицательности и нормировки:

$$\sum_{w \in W_m} \phi_{wt} = 1, \phi_{wt} \geq 0, \sum_{t \in T} \theta_{td} = 1, \theta_{td} \geq 0.$$

EM-алгоритм:

$$p_{tdw} = \mathop{\text{norm}}_{t \in T}(\phi_{wt} \theta_{td});$$

$$\phi_{wt} = \mathop{\text{norm}}_{w \in W} \left(n_{wt} + \phi_{wt} \frac{\partial R}{\partial \phi_{wt}} \right);$$

$$n_{wt} = \sum_{d \in D} \tau_{m(w)} n_{dw} p_{tdw};$$

$$\theta_{td} = \mathop{\text{norm}}_{t \in T} \left(n_{td} + \theta_{td} \frac{\partial R}{\partial \theta_{td}} \right);$$

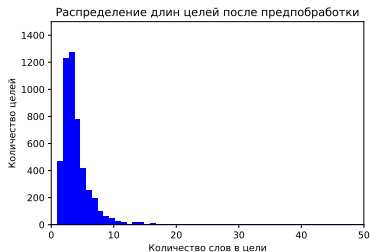
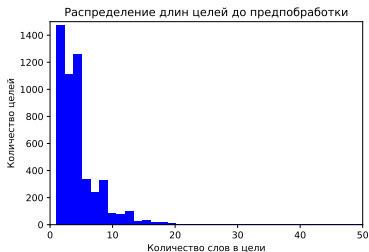
$$n_{td} = \sum_{w \in d} \tau_{m(w)} n_{dw} p_{tdw};$$

В анкетировании приняли участие около 11 тысяч человек. Каждый человек написал по 3 цели. Наиболее подробно люди описывали первую цель. Из полученных 33 тысяч целей было размечено 6 тысяч первых целей.

Данные о целях:

- Формулировка цели на естественном языке
- Уточнение цели (какой может быть первый шаг, конечный результат, преграды для достижения цели, ...)
- Социально-демографическая информация (уровень образования, сфера занятости, пол, возраст, ...)
- Экспертная разметка цели (Specificity, Achievable, Time bound, ...)

Распределение длин целей



	до предобработки	после предобработки
count	5153	4955
min	1	0
50%	3	3
75%	6	4
max	220	139

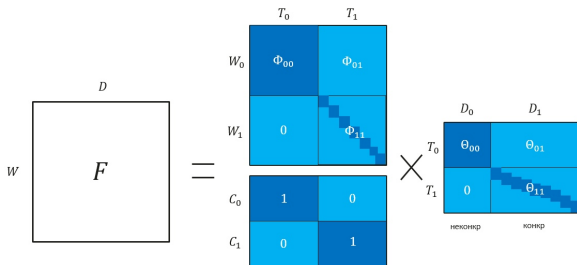
Разработка модели специфичности документов

Имеем коллекцию документов $D = D_0 \cup D_1$.

D_0 - неконкретные документы, D_1 - конкретные, специфичные.

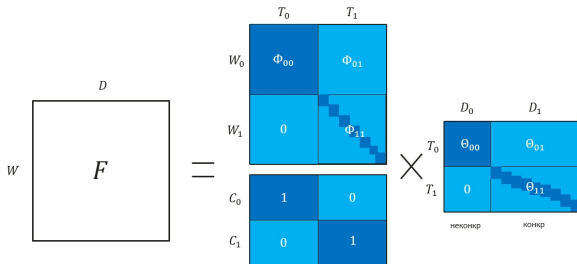
Предположение: Неконкретные слова входят во все документы, конкретные - только в конкретные.

$T = T_0 \cup T_1$, T_0 - неконкретные темы, T_1 - конкретные темы



Инициализация

- Φ_{00} : строим вспомогательную тематическую модель неспецифичных целей, т.е. строим T_0 по D_0
- Φ_{01} : инициализируем нулями
- Φ_{10} , Φ_{11} : инициализируем случайными числами из $U(0, 1)$ так, чтобы вероятности токенов в Φ_{11} были в 100 раз больше, чем в Φ_{01}
- Θ : инициализируем случайными числами из $U(0, 1)$



Обучение

- Φ_{00} : сглаживание
- Φ_{11}, Θ_{11} разреживание
- Φ_{11}, Θ_{11} декорреляция

The diagram shows the decomposition of matrix F into three matrices: W , T , and D . Matrix F is a square matrix with dimension D . Matrix W is a 2x2 block matrix with rows W_0 and W_1 , and columns T_0 and T_1 . Matrix T is a 2x2 block matrix with rows C_0 and C_1 , and columns T_0 and T_1 . Matrix D is a 2x2 block matrix with rows D_0 and D_1 , and columns T_0 and T_1 . The decomposition is shown as $F = W \times T \times D$.

	T_0	T_1
W_0	Φ_{00}	Φ_{01}
W_1	0	Φ_{11}
C_0	1	0
C_1	0	1

\times

	T_0	T_1
D_0	Θ_{00}	Θ_{01}
D_1	0	Θ_{11}
	неконкр	конкр

Параметры, которые можно оптимизировать

- количество неспецифичных тем T_0
- количество специфичных тем T_1
- веса τ_m модальностей
- коэффициенты регуляризации

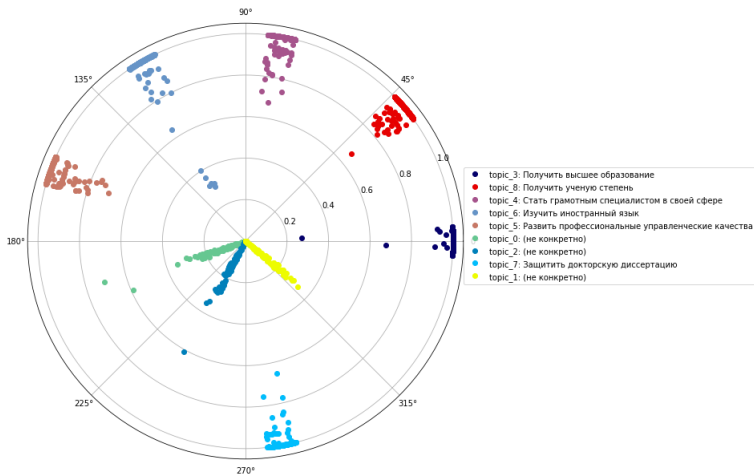
Таблица: Наиболее характерные слова для неконкретных тем

topic_0	знание	topic_1	новое	topic_2	повышение
	навык		статья		уровень
	получение		свой		новый
	развитие		хороший		приобретение
	новый		расширение		умение
	получить		узнать		рост
	новое		кругозор		большой
	компетенция		человек		работа
	новый_знание		узнать_новое		знание_умение
	получение_новый		познание		повышение_уровень

Таблица: Названия конкретных тем

topic_3	Получить высшее образование
topic_4	Стать грамотным специалистом в своей сфере
topic_5	Развить профессиональные управленческие качества
topic_6	Изучить иностранный язык
topic_7	Защитить докторскую диссертацию
topic_8	Получить ученую степень

Рис.: Специфичность всех документов

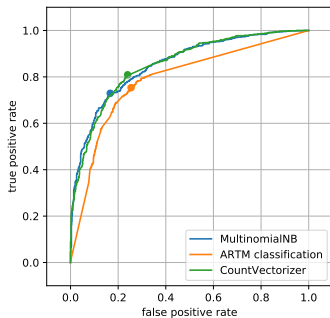


- MultNB — мультиномиальный наивный байесовский классификатор
- LogReg — логистическая регрессия
- ARTM_class — простая ARTM для классификации
- Main_model — предлагаемая ARTM

критерий	Main_model	MultNB	LogReg	ARTM_class
roc_auc	0.871657	0.843594	0.843779	0.787452
precision	0.752055	0.788851	0.792321	0.705634
recall	0.825564	0.702256	0.682707	0.753383
accuracy	0.800269	0.782784	0.778077	0.749159

- 1 Вероятностная тематическая модель мультимодальных анкетных данных с частичным обучением для оценивания тематики и конкретности образовательных целей пользователей в системе дистанционного образования.
- 2 Способ обучения тематической модели, предполагающий явное разбиение коллекции документов, множества тем и словаря слов на специфичные и неконкретные.
- 3 Способ графической визуализации тематической модели в виде круговой диаграммы «темы–специфичность».
- 4 Результаты экспериментов, показывающие, что предложенная модель классифицирует анкетные данные точнее, чем стандартные модели бинарной классификации.

Спасибо за внимание!



	MultNB	LogReg	ARTM_classification
roc_auc	0.843594	0.843779	0.787452
precision	0.788851	0.792321	0.705634
recall	0.702256	0.682707	0.753383
accuracy	0.782784	0.778077	0.749159

Подбор параметров модели

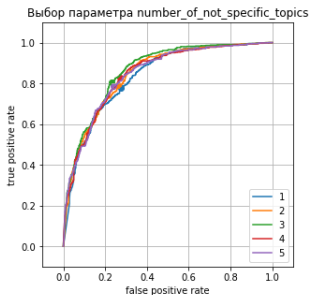


Рис.: Зависимость качества модели от количества неспецифичных тем

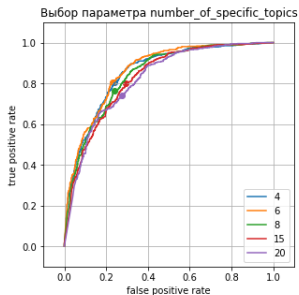


Рис.: Зависимость качества модели от количества специфичных тем