

# 京东-贪心 NLP 项目实验手册

## 项目 3：京东客服对话系统

项目设计和编写：林培鑫 姜冰钰

后期编辑：李文哲

单 位：贪心科技

2020 年 8 月 27 日

# 目 录

<b>1</b>	<b>项目描述与目标</b>	<b>3</b>
<b>2</b>	<b>项目数据描述</b>	<b>4</b>
<b>3</b>	<b>项目事宜</b>	<b>4</b>
3.1	项目的整个框架 . . . . .	5
<b>4</b>	<b>项目主要技术介绍</b>	<b>7</b>
4.1	Average word vectors . . . . .	7
4.2	Hierarchical Navigable Small World . . . . .	8
4.3	Learning-to-Rank . . . . .	9
4.4	BERT . . . . .	10

贪心科技版权所有

# 1 项目描述与目标

**对话系统 (Dialogue Systems)** 又可以称之为**聊天机器人 (ChatBot)**，主要是实现自动与用户进行对话的功能。帮助用户完成某些具体的任务（下单、打车、订座等）的对话系统可以称之为**任务导向型 (Task-oriented)** 的对话系统；解答用户的某些问题（**询问天气、股价、交通等**）的对话系统可以称之为**问答型 (QA-based)** 的对话系统；除此之外，还有和用户聊天的**闲聊型 (Chatting)** 对话系统。大多数的对话系统都是混合了几种类型的功能。

对话系统中语言的生成主要可以分为两种方式：**检索式 (Retrieval)** 和**生成式 (Generative)**。**检索式方法**我们一般会构建一个语料库，为 FAQ 存放 query-response pairs，然后用户发起一个新的 query 时，我们去匹配为这个 query 检索最佳的 response。上述过程一般又可以分为召回 (Retrieve) 和排序 (ranking) 两个部分：召回即通过 query 找到语料库中最相似的几十个或几百个 query，大大减小候选 Response 的数量。这一部分我们一般采用轻量级的方法，如倒排索引 (Inverted Index) 和近似近邻搜索 (Approximate Nearest Neighbor Search) 等进行快速检索。排序则是对召回的结果进行进一步的筛选，可以构建更复杂的特征，使用机器学习或深度学习的方法来进行排序。**生成式方法**基本可以使用我们在上一阶段的文本生成项目所用的 Seq2seq 或 PGN 技术，但是现在更热门的是基于 Self-Attention 的预训练模型 (BERT、GPT2 等)。另外，对话中回答的生成会比摘要或者翻译的难度大很多，原因在于给定 Encoder 的输入，在这一任务中 Decoder 的候选序列比起摘要或翻译要多得多。[1]

通过本项目的练习，你能掌握对话系统的技术：

- **意图识别**：这个地方我们做成一个简单的文本二分类任务，根据用户的开场白识别用户的意图是业务需求还是闲聊。
- **检索模型**：这个项目中我们用检索式的方法来做业务型对话的回答生成。我们将用 Approximate Nearest Neighbor Search(ANNS) 方法中较为常用的 Hierarchical Navigable Small World(HNSW) 来做召回的部分；然后将构建各种相似度特征（包括深度匹配网络），并利用

LightGBM 来训练一个 Learning2Rank 模型。

- **生成模型**：这个项目中我们将使用 BERT 模型来实现闲聊对话的生成。然后对于体积过大的 BERT 模型，我们还会进行压缩以便进行部署。
- **对话管理**：对于多轮的对话，我们需要对考虑上一轮的内容，以便对下一轮生成的回复进行合理性的控制。

同时，通过本项目，你将：

- 1. 熟练掌握如何使用 FastText 做文本分类。
- 2. 熟练掌握 ANNS 理论和技术。
- 3. 熟练掌握 Learning2Rank 的特征工程和 LightGBM 模型。
- 4. 熟练掌握 Transformer 和 BERT 模型。
- 5. 熟练掌握模型压缩的方法。

作为第三个项目，我们的目标是帮助大家完成一个客服智能对话系统，通过这个任务让大家熟悉对话系统的基本架构和常用技术。该项目的数据和部分技术用于京东的 JIMI 和京小智。

## 2 项目数据描述

**京东对话数据集**在本项目中，我们使用的是京东电商的客服对话数据。原始数据有将近 200 万行，如图 1所示。

我们会只取大概 10% 的数据，对同属一个对话的句子进行合并，两两结合为 query-response pairs，如图 2所示。

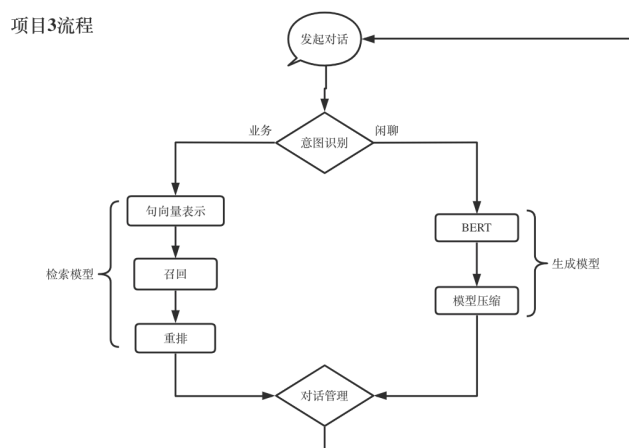
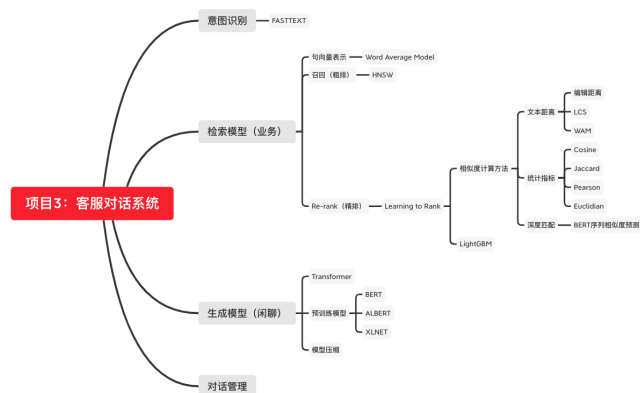
## 3 项目事宜

本项目相对于上一个项目，将更加注重对大家动手能力的提升，让大家完整的去了解一个 NLP 项目的整体框架，并且实现各个环节的重要细节。

[illegible][illegible]

### 3.1 项目的整个框架

- **意图识别**: 训练一个 FastText 模型对用户的输入进行意图识别, 对业



务型需求使用检索模型来生成回复，对闲聊需求则使用生成模型来生成回复。这一部分的训练数据我们使用关键词来自动标注，每个样本是一句用户的输入，整理一系列业务相关的词汇，将包含业务关键词的样本标注为 1，否则为 0。

- **召回**：我们使用 Word Average Model (WAM) 来表示句向量，之后用 HNSW 模型来实现 ANNS。HNSW 我们可以使用 hnsplib 或者 faiss 来实现。

- **排序**：这一步我们需要构建多种相似度特征，主要可以分为几类：1. 基于字符串距离的（编辑距离、列文斯坦距离、LCS）；2. 基于向量距离的（cosine、Euclidian、Jaccard、WMD）；3. 基于统计量的（BM25、Pearson Correlation）；4. 基于深度匹配模型的。构建完特征后，我们使用 LightGBM 来训练一个 LTR 模型。
- **生成模型**：这个地方我们使用 BERT（也可以是 GPT2 等预训练模型）来做 seq2seq 的生成任务。之后对 BERT 模型进行压缩（蒸馏或剪裁）。
- **对话管理**：对轮对话管理需要实现 Dialogue State Tracing (DST) 和 Policy Selection。

## 4 项目主要技术介绍

### 4.1 Average word vectors

为了通过计算距离来实现近似搜索，我们需要将一个句子或者一篇文档表示为一个向量，这种方式叫做 Sentence Embedding 或者 Doc Embedding。其中最常用的一种模型是 Average word vectors 或者叫 Word Averaging Model(WAM)。其方法是得到将一个句子中每一个词向量在同一个 embedding 维度的值取平均值得到该句子在该嵌入维度的值。如图 5所示。

$$\begin{array}{c} W_1 \\ \begin{bmatrix} W_{11} \\ W_{12} \\ \vdots \\ W_{1n} \end{bmatrix} \end{array} + \begin{array}{c} W_2 \\ \begin{bmatrix} W_{21} \\ W_{22} \\ \vdots \\ W_{2n} \end{bmatrix} \end{array} + \dots + \begin{array}{c} W_n \\ \begin{bmatrix} W_{n1} \\ W_{n2} \\ \vdots \\ W_{nn} \end{bmatrix} \end{array} = \begin{array}{c} D \\ \begin{bmatrix} \frac{W_{11}+W_{21}+\dots+W_{n1}}{n} \\ \vdots \\ \frac{W_{1n}+W_{2n}+\dots+W_{nn}}{n} \end{bmatrix} \end{array}$$

图 5: Average word vectors

## 4.2 Hierarchical Navigable Small World

HNSW[4] 的前身是 NSW (Navigable-Small-World-Graph)。NSW 通过设计出一个具有导航性的图来解决近邻图发散搜索的问题，但其搜索复杂度仍然过高，达到多重对数的级别，并且整体性能极易被图的大小所影响。HNSW 则是着力于解决这个问题。作者借鉴了 SkipList 的思想，提出了 Hierarchical-NSW 的设想。简单来说，按照一定的规则把一张的图分成多张，越接近上层的图，平均度数越低，节点之间的距离越远；越接近下层的图平均度数越高，节点之间的距离也就越近。搜索从最上层开始，找到本层距离最近的节点之后进入下一层。下一层搜索的起始节点即是上一层的最近节点，往复循环，直至找到结果（见图 6）。由于越是上层的图，节点越是稀少，平均度数也低，距离也远，所以可以通过非常小的代价提供了良好的搜索方向，通过这种方式减少大量没有价值的计算，减少了搜索算法复杂度。更进一步，如果把 HNSW 中节点的最大度数设为常数，这样可以获得一张搜索复杂度仅为  $\log(n)$  的图。

实现 HNSW 主要有以下两个 package 可以选用：

**hnswlib**：一个轻量级的 library，专门用来实现 HNSW 算法，底层用 C++ 实现，有实现 Python 的 wrapper。

**Facebook AI Similarity Search (Faiss)**：Facebook 开发的专门用于做 Proximate Search 的 library，支持包括 KNN 和 HNSW 等多种搜索算法，也支持 PCA 等降维算法以提升性能。而且支持 GPU。Faiss 的底层也是用 C++ 开发，并实现 Python wrapper。Github repo: <https://github.com/facebookresearch/faiss>。



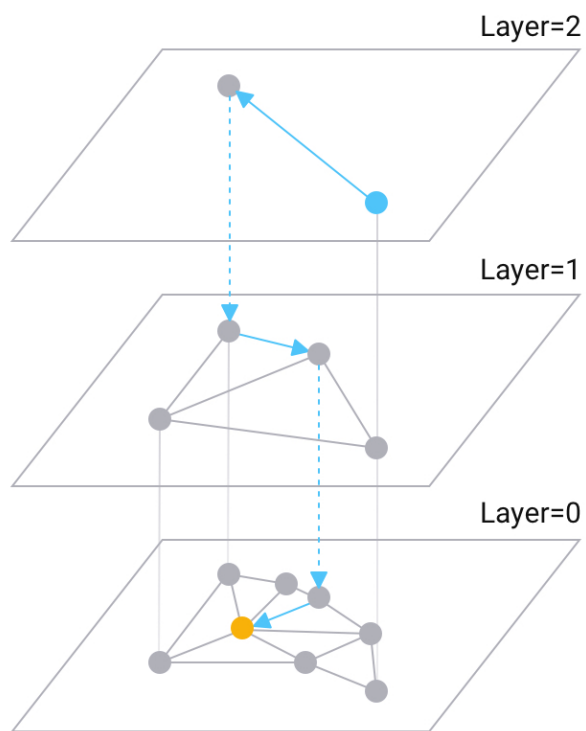


图 6: Hierarchical Navigable Small World

### 4.3 Learning-to-Rank

排序学习 [3] 是一个有监督的机器学习过程，对每一个给定的查询 - 文档对，抽取特征，通过日志挖掘或者人工标注的方法获得真实数据标注。然后通过排序模型，使得输入能够和实际的数据相似。常用的排序学习分为三种类型：PointWise, PairWise 和 ListWise。如图 5所示。

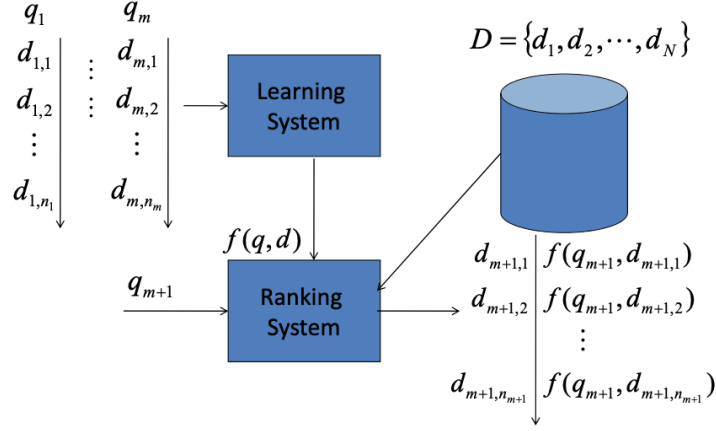


图 7: Learning-to-Rank

#### 4.4 BERT

BERT[2] 是 Google 于 2018 年提出的语言模型，在多种 NLP 任务上有出色的表现，因此近几年成为 NLP 领域的热门模型。BERT 的全称是 Bidirectional Encoder Representation from Transformers，模型的核心有两个：1 是使用了基于 Self-Attention 的 Transformer[5]；2 是使用了 pre-train 的方法，即用了 Masked LM 和 Next Sentence Prediction 两种方法分别捕捉词语和句子级别的 representation。如图 8，9 所示。

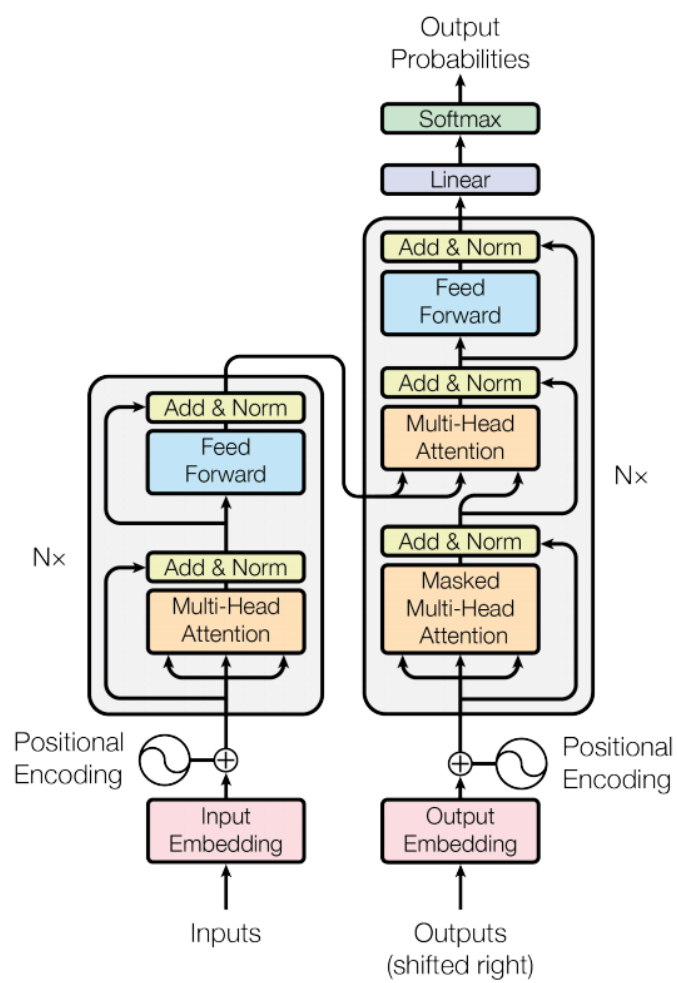


图 8: Transformer

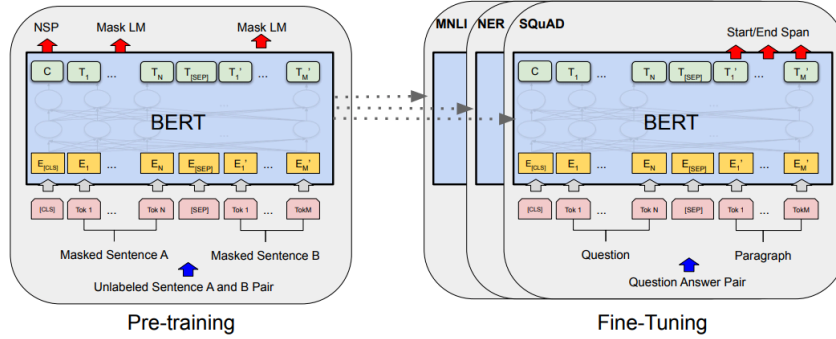


图 9: Bert

## 参考文献

- [1] CHEN, H., LIU, X., YIN, D., AND TANG, J. A survey on dialogue systems: Recent advances and new frontiers. *CoRR abs/1711.01731* (2017).
- [2] DEVLIN, J., CHANG, M., LEE, K., AND TOUTANOVA, K. BERT: pre-training of deep bidirectional transformers for language understanding. *CoRR abs/1810.04805* (2018).
- [3] LI, H. A short introduction to learning to rank. *IEICE Trans. Inf. Syst.* 94-D (2011), 1854–1862.
- [4] MALKOV, Y. A., AND YASHUNIN, D. A. Efficient and robust approximate nearest neighbor search using hierarchical navigable small world graphs. *CoRR abs/1603.09320* (2016).
- [5] VASWANI, A., SHAZEER, N., PARMAR, N., USZKOREIT, J., JONES, L., GOMEZ, A. N., KAISER, L., AND POLOSUKHIN, I. Attention is all you need. *CoRR abs/1706.03762* (2017).