

OPEN DATA SCIENCE CONFERENCE

San Francisco | Oct. 28 - Nov. 1, 2019



Building an Industry Classifier

Scrapy / BERT / PyTorch / SageMaker

Ido Shlomo - BlueVine

About BlueVine

- Fintech startup up based in Redwood City, CA
- Provides working capital (loans) to small & medium sized businesses
- Over \$2 BN funded to date
- Mission of DS team:
 - Automate
 - Scale
 - Improve

The image shows a screenshot of the BlueVine website. At the top right is the BlueVine logo and a three-line menu icon. Below the header, the text "Fast funding for your business." is displayed in bold. A blue button labeled "Apply now" is positioned below the text. A small note states "Applying will not affect your credit score¹". Below this is a photograph of a smiling man wearing a light-colored shirt and dark overalls. At the bottom, a blue banner contains the text "Business Loans from \$5,000 - \$5 million".

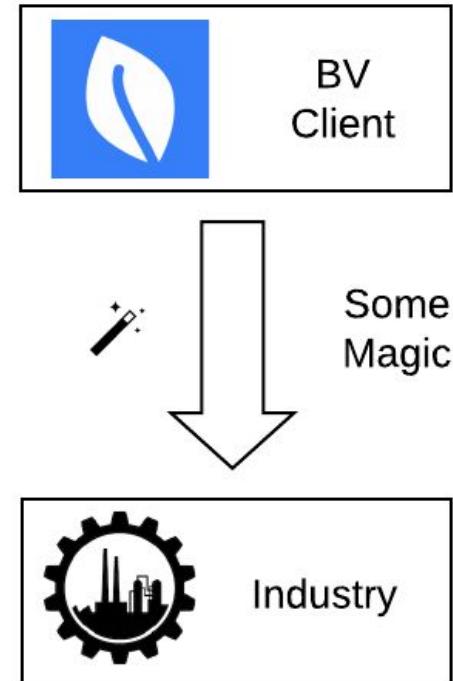
About Me

- Data Science Manager @ BlueVine
- Lead BlueVine's DS team in Redwood City, CA (total of ~20 people across RWC & TLV)
- Team focus:
 - Risk / Fraud
 - Financial Strength
 - Marketing
- Personal interests: Unstructured data and DS Infrastructure.

Slides at: https://github.com/ido-sh/odsc_presentations

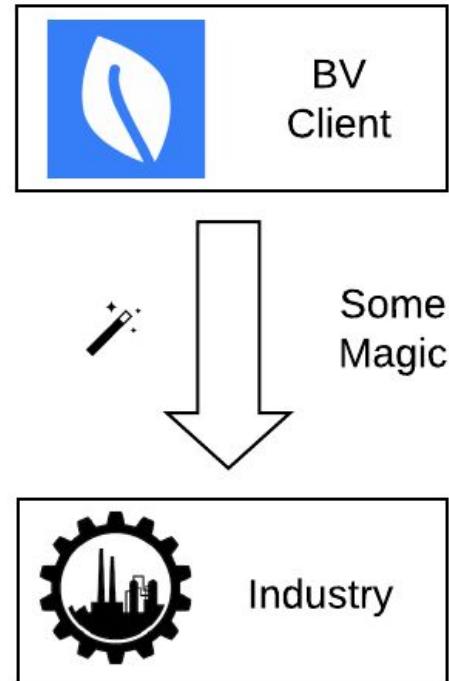
The Point of this Presentation

- Business need: For every client → Predict their Industry
- Industry can be one of 135 categories, e.g:
 - Electrician
 - Marketing
 - Steel manufacturing
- Problem: Traditional sources are pricey, inaccurate and not widely available

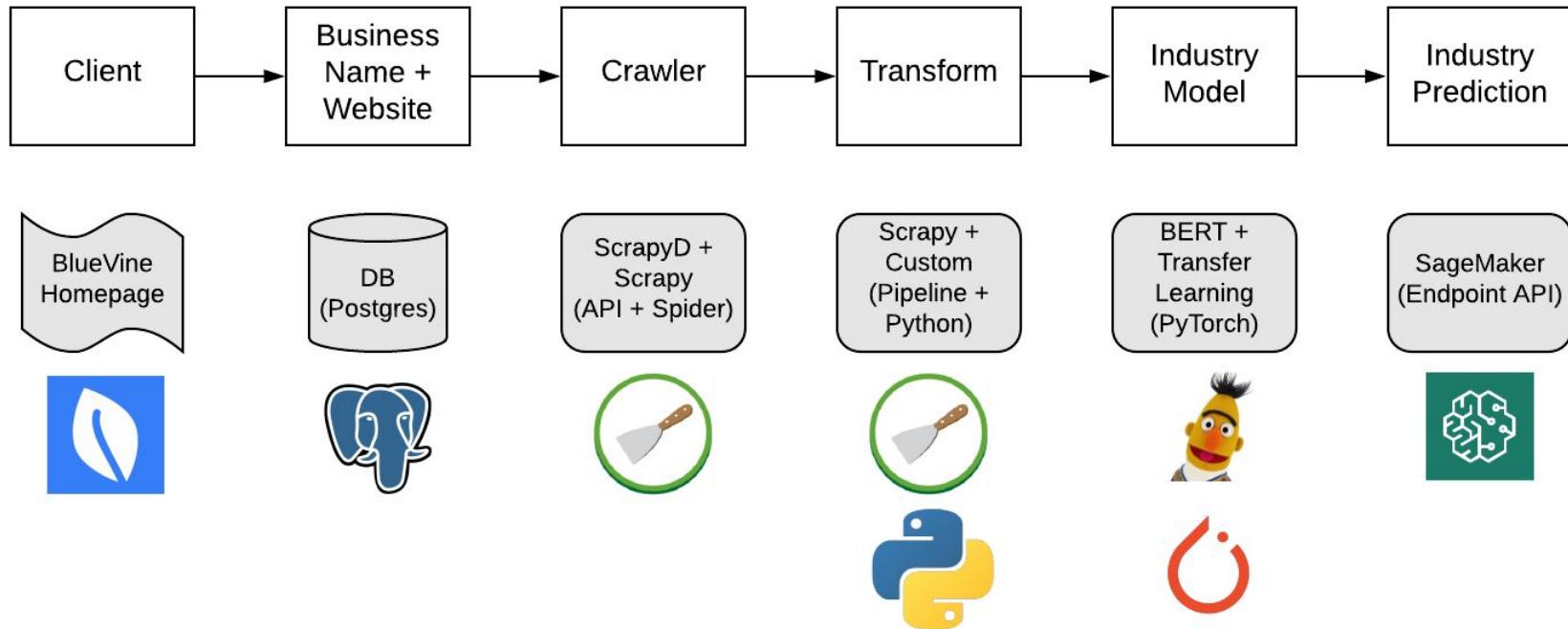


The Point of this Presentation

- Our solution: Infer industry from the business website
- Main Challenges
 - High quality sample
 - Efficient crawling
 - Strong NLP model
 - Scalable & independent deployment
- Main tools: Scrapy D, Scrapy, PyTorch, BERT, Amazon SageMaker



Inference Pipeline



Run example

<end to end notebook>

Presentation Overview

1. Target

- Gathering websites
- Tagging the industry

3. Model

- Design (BERT)
- Train
- Deploy

2. Data

- Building the crawler
- Crawling at scale

4. Conclusion

- Summary
- Key takeaways

Target - Gathering websites

Potential sources:

- Internal BV data
 - E.g: Client provided information
 - Not very big
 - Skewed scope
- Proprietary Databases
 - E.g: Various vendors
 - Big
 - Pricy
 - Outdated
- Directed web searches
 - E.g: “car dealerships in the US”
 - High effort
 - Limited scope

Target - Gathering websites

Problem:

- We have 135 industries
- Want sample to be balanced
- In order to balance, need to know the Industry
- If we knew the Industry, wouldn't need the model...

Target - Gathering websites

Problem:

- We have 135 industries
- Want sample to be balanced
- In order to balance, need to know the Industry
- If we knew the Industry, wouldn't need the model...

Solution:

- Use proxy to “guess” the industry
 - Internal BV Data: User provided
 - Web Searches: Search term
 - Proprietary Data: NAICS codes
- Fingers crossed, hope for the best...

Target - Tagging the Industry

Goals:

- Get quality tags
- Reasonable Timeframe
- Stay within budget

Crowdsourcing challenges:

- Understanding 135 industry categories is hard
- Understanding what a business does is hard
- Tie-breaking is hard

Row Id	Url	Biz Name	Industry
1	a.com	a inc	X
2	b.com	b corp	Y
3	c.bom	c llc	Z
4	d.com	d	X
5	e.com	e	Y
...

Target - Tagging the Industry

Crowdsourcing Strategies:

- Break up the task into smaller parts:
 - One obs requires multiple decisions
 - More tasks → more \$\$\$
- Train team for complex decision making:
 - Time intensive → more \$\$\$
 - Slow
- Get multiple judgements for every website:
 - If everybody are bad, result will still be bad
 - More judgements → more \$\$\$

Target - Tagging the Industry

What we did:

- Set max \$\$\$ Budget
- Set max time per task (in seconds)
- Developed small in-house tagged test set
- Held vendor trials for all strategies
- From vendors within max time cap: chose one with highest accuracy
- PLUG: CloudFactory (they are great)

<https://www.cloudfactory.com/>

Target - Tagging the Industry

What we did:

- Set max \$\$\$ Budget
- Set max time per task (in seconds)
- Developed small in-house tagged test set
- Held vendor trials for all strategies
- From vendors within max time cap: chose one with highest accuracy
- PLUG: CloudFactory (they are great)

<https://www.cloudfactory.com/>

Outcome:

- Overall time: 3 months
- Overall size: 250K observations (business names + websites)
- NOT balanced, BUT sufficient
- Fingers crossed, hope for the best...

Data - Building the Crawler

For each URL, we need:

- Home page text cleaned (no HTML)
- Internal page(s) text cleaned (no HTML)
- Page metadata
 - E.g: Path, title, click text, etc.
- Ability to choose
 - E.g. do I even want this page?
- Raw data as backup

Data - Building the Crawler

For each URL, we need:

- Home page text cleaned (no HTML)
- Internal page(s) text cleaned (no HTML)
- Page metadata
 - E.g: Path, title, click text, etc.
- Ability to choose
 - E.g. do I even want this page?
- Raw data as backup

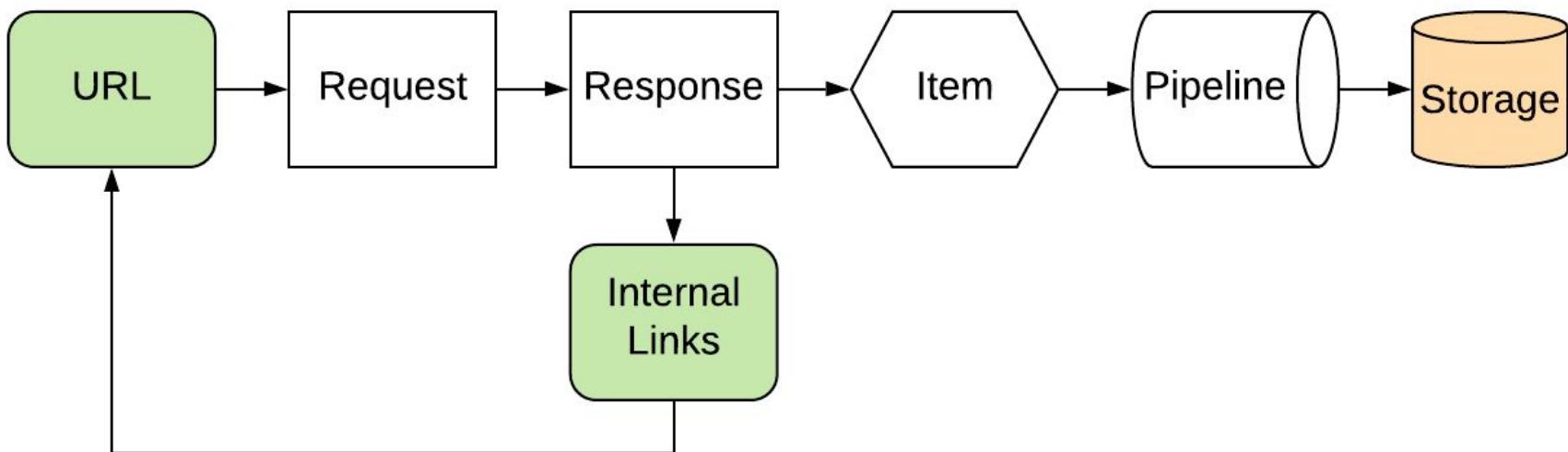
Our tool:

- Scrapy: fast high-level web crawling and web scraping framework

<https://github.com/scrapy/scrapy>

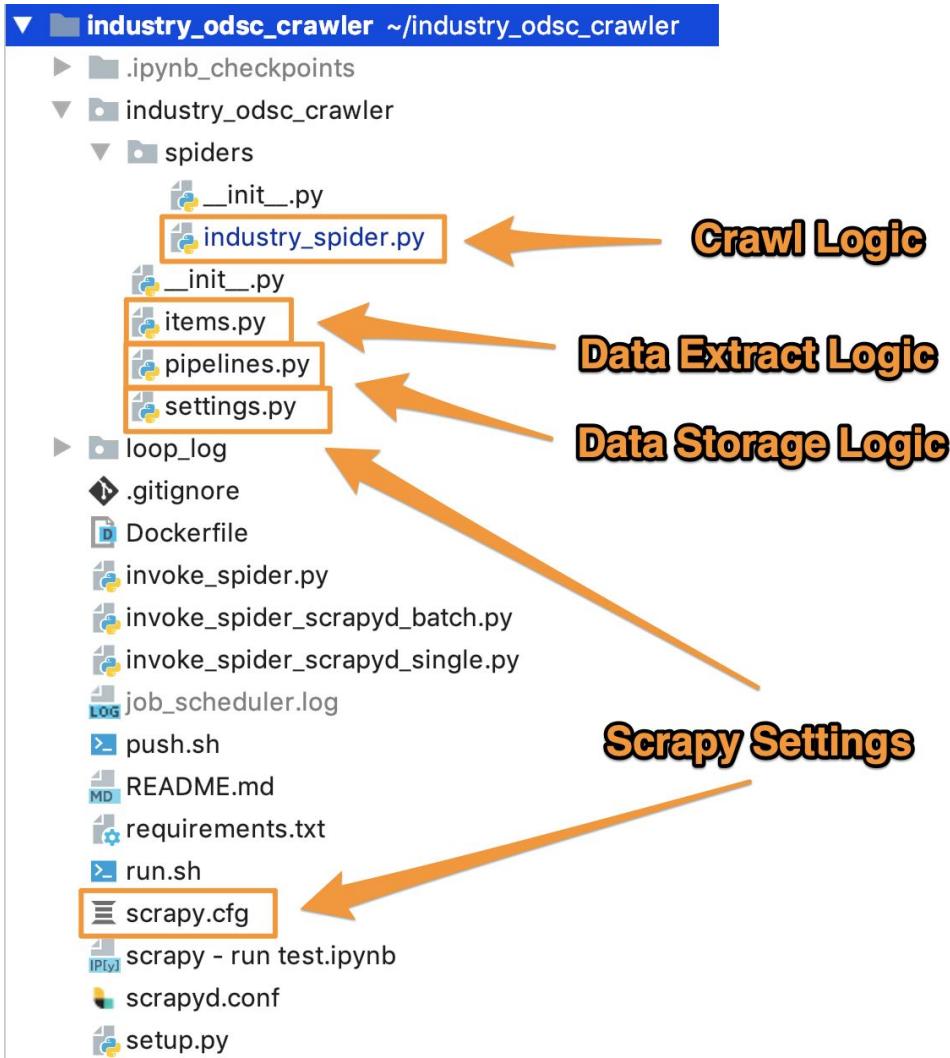
Data - Building the Crawler

Scrapy Flow: Single URL



Data - Building the Crawler

<scrapy - run test notebook>



```
14 class IndustrySpider(Spider):
15     name = 'industry_v4_spider'
16
17     def __init__(self, ids, urls, save_local=False,
18                  **kwargs):
19         self.start_urls = urls.split(",")
20         self.ids = [int(x) for x in ids.split(",")]
21         self.link_counter = dict(zip(self.ids, [0] * len(self.ids)))
22         self.save_local = save_local == 'True'
23         super().__init__(**kwargs)
24
25     def start_requests(self):
26         for url_id, url in zip(self.ids, self.start_urls):
27             try:
28                 request = Request(url=url, callback=self.parse,
29                                    dont_filter=True,
30                                    errback=self.handle_error)
31             except ValueError:
32                 request.meta.update(url_id=url_id, url_order=1,
33                                     url_rank=1,
34                                     url_click_text=None,
35                                     url_domain=urlparse(url).netloc)
36
37             yield request
```

URLs to start with

Internal link counter

Make 1st request per URL

```
42 def parse(self, response):
43     web_page = WebPage()
44     web_page['url_id'] = response.meta['url_id']
45     web_page['url_order'] = response.meta['url_order']
46     web_page['input_url'] = response.meta['download_slot']
47     web_page['depth'] = response.meta['depth'] if \
|       'depth' in response.meta.keys() else None
48     web_page['url'] = response.url
49     web_page['path'] = urlparse(response.url).path
50     web_page['url_click_text'] = response.meta['url_click_text']
51     web_page['response_raw'] = response.text
52     web_page['url_rank'] = response.meta['url_rank']
53     web_page['response_text'], web_page['url_title'] = \
|       text_from_html(response.text)
54     web_page['status_code'] = response.status
55     web_page['is_error'] = False
56     self.link_counter[response.meta['url_id']] += 1
57     yield web_page
58
59
```



Extract only text

```
61 allowed_domains = {response.meta['url_domain']}
62 if 'redirect_urls' in response.meta.keys() \
63     and isinstance(response.meta['redirect_urls'], list):...
69
70     allowed_domains = {x.replace("www.", "") for x in allowed_domains}
71     links = LinkExtractor(canonicalize=True, unique=True).\
72         extract_links(response)
73     links = [link for link in links if self.is_link_allowed(
74         allowed_domains, link.url)]
75     links = sorted(links, key=self.get_link_rank) ← Sort by "priority"
76
77 for link in links:
78     is_below_limit = self.link_counter[response.meta['url_id']] < \
79         MAX_PAGES_PER_DOMAIN
80     if is_below_limit:... ↑ Stay ⇌ X links per domain
```

Extract Internal
Links

Sort by "priority"

Stay ⇌ X links per domain

```
89  
90 @staticmethod  
91     def get_link_rank(link):  
92         good_hits = [  
93             ('about', 2),  
94             ('story', 2),  
95             ('what we do', 3),  
96             ('capabilities', 3),  
97             ('services', 4),  
98             ('products', 4),  
99             ('amenities', 5),  
100            ('overview', 6)  
101        ]  
102        url_rank = 100  
103        text_rank = 100  
104        for hit, rank in good_hits:  
105  
106  
107        for hit, rank in good_hits:  
108  
109  
110  
111        link_rank = min(url_rank, text_rank)  
112        return link_rank
```

Prioritize
informative
pages

Default value
(lower is better)

```
9 class ConvertToDataFramePipeline:  
10    def __init__(self):  
11        self.item_list = []  
12  
13    def process_item(self, item, spider):  
14        self.item_list += [item]  
15  
16    def close_spider(self, spider):  
17        data_clean_json, fn_clean = combine_items(  
18            self.item_list, spider, is_raw=False)  
19  
20        data_raw_json, fn_raw = combine_items(  
21            self.item_list, spider, is_raw=True)  
22  
23    if spider.save_local: ...  
24    else:  
25        fp = 'test_crawls/{}'.format(fn_clean)  
26        fp_raw = 'test_crawls/{}'.format(fn_raw)  
27        bucket_name = 'industry-odsc-19'  
28        s3 = boto3.resource('s3')  
29        bucket = s3.Bucket(bucket_name)  
30        logger.info("saving to s3: {}/{}".format(bucket_name, fp))  
31        data_clean_json_obj = dump_json_to_object(data_clean_json)  
32        bucket.put_object(Body=data_clean_json_obj, Key=fp)  
33        logger.info("saving to s3: {}/{}".format(bucket_name, fp_raw))  
34        data_raw_json_obj = dump_json_to_object(data_raw_json)  
35        bucket.put_object(Body=data_raw_json_obj, Key=fp_raw)  
36  
37  
38  
39  
40  
41
```

Collect all items from spider

Process "clean" data for model

Keep "raw" as backup

Save in S3

Data - Crawling at Scale

For our list of 250K URLs, we need to:

- Run batches, in parallel
- Persistent process that's always up
- Some logging / UI mechanism to follow progress

Data - Crawling at Scale

For our list of 250K URLs, we need to:

- Run batches, in parallel
- Persistent process that's always up
- Some logging / UI mechanism to follow progress

Our tool:

- ScrapyD: a service for running Scrapy spiders
<https://github.com/scrapy/scrapyd>
- Deploy Scrapy projects and control their spiders using an HTTP JSON API

```
1 [scrapyd]
2 eggs_dir      = eggs
3 logs_dir      = logs
4 items_dir     =
5 jobs_to_keep  = 500
6 dbs_dir       = dbs
7 max_proc      = 0
8 max_proc_per_cpu = 4
9 finished_to_keep = 100
10 poll_interval = 5.0
11 bind_address  = 0.0.0.0
12 http_port    = 6800
13 debug        = off
14 runner       = scrapyd.runner
15 application  = scrapyd.app.application
16 launcher     = scrapyd.launcher.Launcher
17 webroot      = scrapyd.website.Root
18
19 [services]
20 schedule.json = scrapyd.webservice.Schedule
21 cancel.json   = scrapyd.webservice.Cancel
22 addversion.json = scrapyd.webservice.AddVersion
23 listprojects.json = scrapyd.webservice.ListProjects
24 listversions.json = scrapyd.webservice.ListVersions
25 listspiders.json = scrapyd.webservice.ListSpiders
26 delproject.json = scrapyd.webservice.DeleteProject
27 delversion.json = scrapyd.webservice.DeleteVersion
28 listjobs.json = scrapyd.webservice.ListJobs
29 daemonstatus.json = scrapyd.webservice.DaemonStatus
```

Parallelism settings

Webserver port

```
1 #!/usr/bin/env bash  
2 set -m  
3 scrapyd > /dev/null 2>&1 &  
4 sleep 5  
5 scrapyd-deploy  
6 fg %1
```

Start Server

**Deploy Project
(Spider)**

```
--  
19      from scrapyd_api import ScrapydAPI  
20  
21  def chunks(l, n):  
22      """Yield successive n-sized chunks from l."""  
23      for i in range(0, len(l), n):  
24          yield l[i:i + n]  
25  
26  
27  scrapyd_host = 'http://localhost:6800'  
28  project_name = 'industry_odsc_crawler'  
29  spider_name = 'industry_v4_spider'  
30  scrapyd = ScrapydAPI(scrapyd_host)  
31  
32  data_fn = 'big_url_list.csv'  
33  data_fp = os.path.join(os.path.dirname(__file__),  
34                      'data', data_fn)  
35  data_df = pd.read_csv(data_fp)  
36  
37  url_id_chunks = chunks(data_df['url_id'], 100)  
38  url_id_chunks = list(enumerate(url_id_chunks))  
39  total_num_chunks = len(url_id_chunks)
```

Import ScrapyD API constructor

Specify ScrapyD host and Scrapy spider

Load big list of URLs

Split into chunks of 100

```
41 max_running = 10          ← Max # concurrent batches
42 max_pending = 5           ← Max # batches in queue
43 sleep_interval = 15       ← Server polling rate
44
45 while len(url_id_chunks) > 0:
46     server_status = scrapyd.list_jobs(project_name)   ← Check which jobs,
47                                         running, pending, etc
48
49     num_pending = len(server_status['pending'])
50     num_running = len(server_status['running'])
51
52     if (num_pending > max_pending) or (num_running > max_running):... ← If bounds
53                                         exceeded, wait
54
55     i, url_id_chunk = url_id_chunks.pop(0)
56     urls = data_df.loc[data_df['url_id'].isin(url_id_chunk), 'url'].tolist()
57
58     ids = ','.join([str(x) for x in url_id_chunk])
59     urls = ','.join(urls)
60
61     res = scrapyd.schedule(project_name, spider_name,
62                           ids=ids, urls=urls)   ← If slot open,
63                                         schedule a batch
64
65     logger.info("added job {} (chunk num: {} / {})".format(
66         res, i, total_num_chunks))
```

Model - Design (BERT)

BERT (grossly simplified):

- Can turn string of text into a vector of length 768 (or 1024)
- Max 512 tokens allowed
- Open source pre-trained models available for transfer learning (cased, uncased, base, large)
- For PyTorch:

<https://github.com/huggingface/transformers>

Model - Design (BERT)

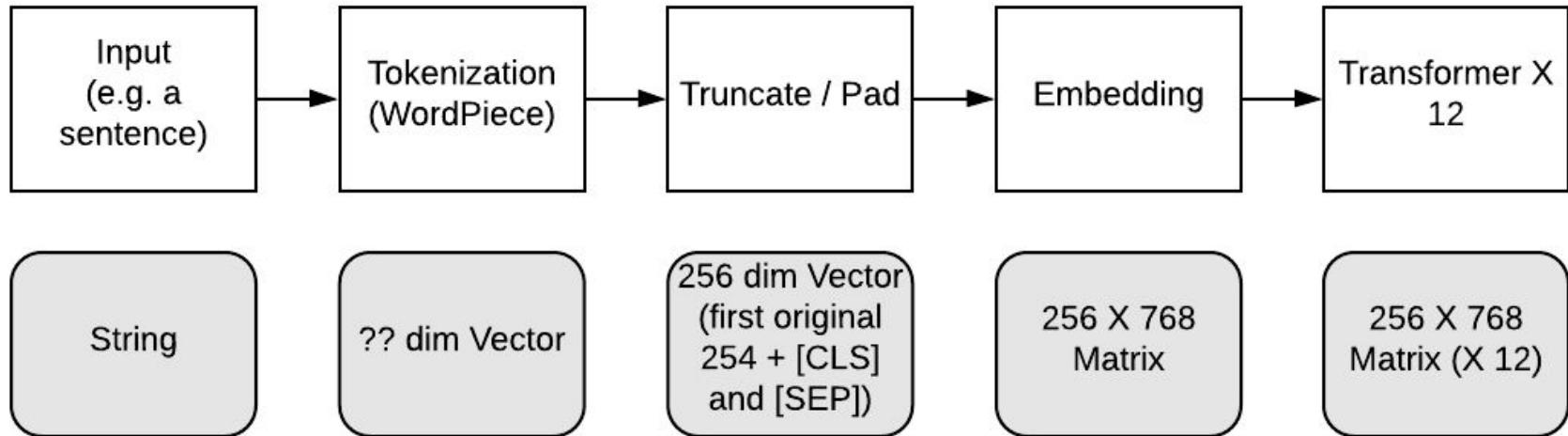
BERT (grossly simplified):

- Can turn string of text into a vector of length 768 (or 1024)
- Max 512 tokens allowed
- Open source pre-trained models available for transfer learning (cased, uncased, base, large)
- For PyTorch:
<https://github.com/huggingface/transformers>

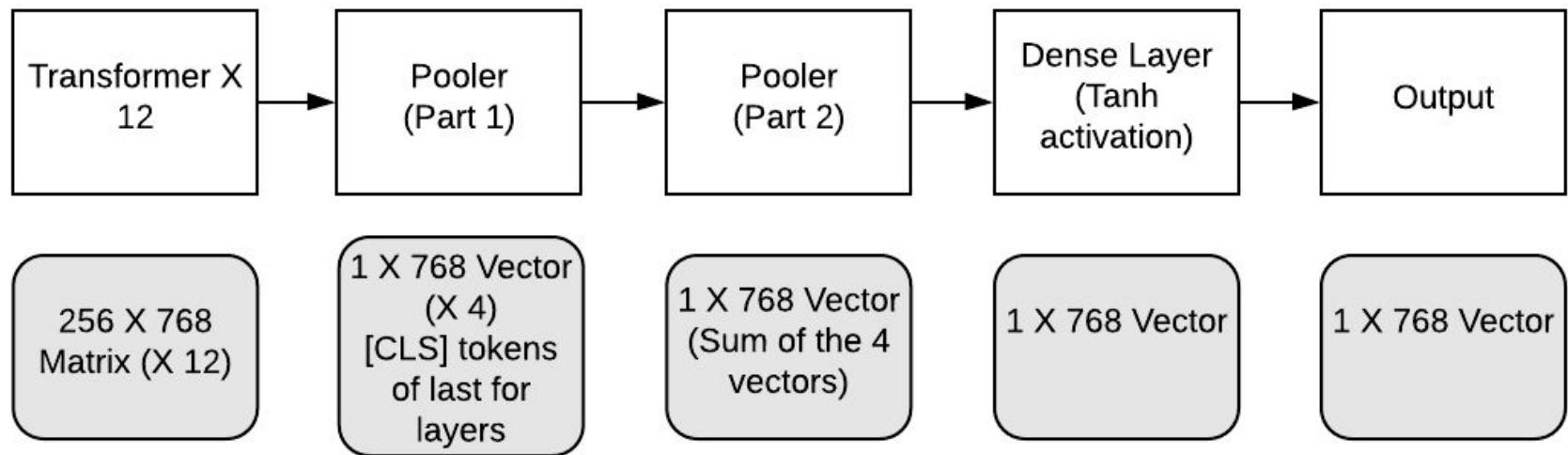
A bit more context:

- Uses WordPiece tokenization to split words into tokens
- Uses Transformers to perform “attention”, i.e: consider same token in different contexts
- Bi directional
- Achieves SOTA performance on benchmark NLP tasks

Model - Design (BERT)



Model - Design (BERT)



Model - Design

The basic idea:

- A single observation is a pair of:
Business name & Website
- Organize data for every obs into a clear
set of inputs
- Use BERT to convert all free text into
numeric vectors
- Wrap this into a single model to classify
our 135 industries

Model - Design

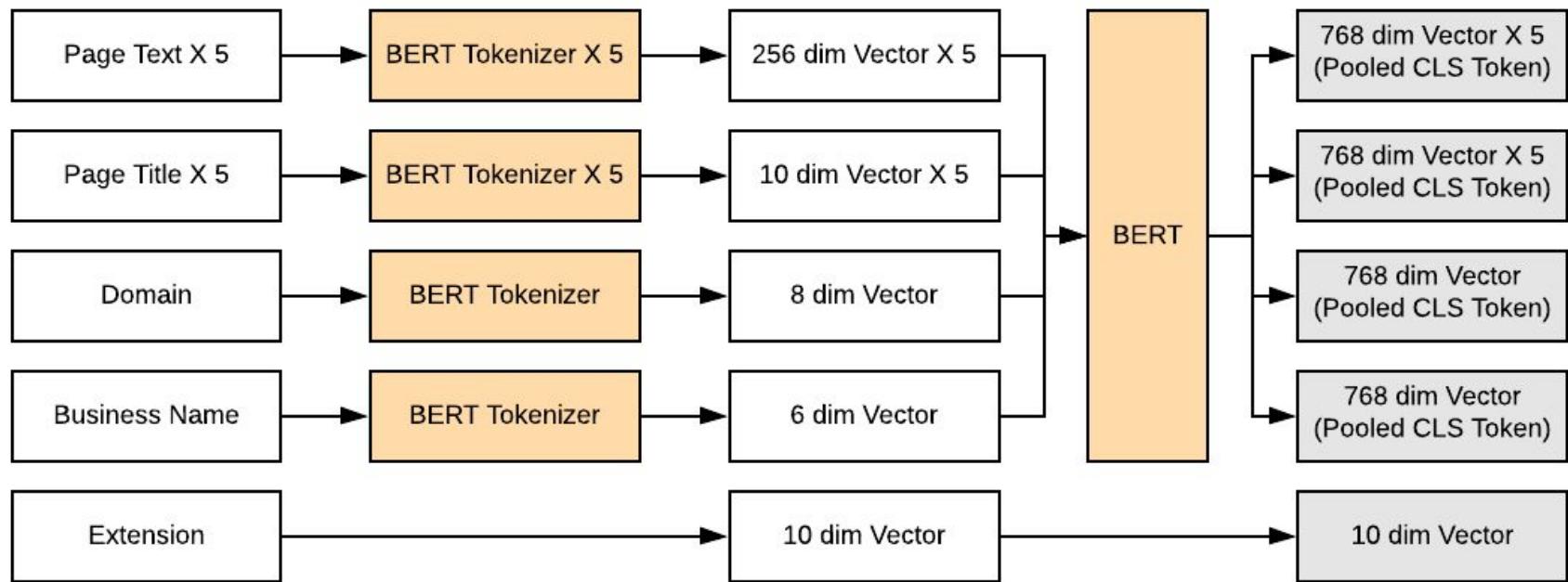
The basic idea:

- A single observation is a pair of:
Business name & Website
- Organize data for every obs into a clear set of inputs
- Use BERT to convert all free text into numeric vectors
- Wrap this into a single model to classify our 135 industries

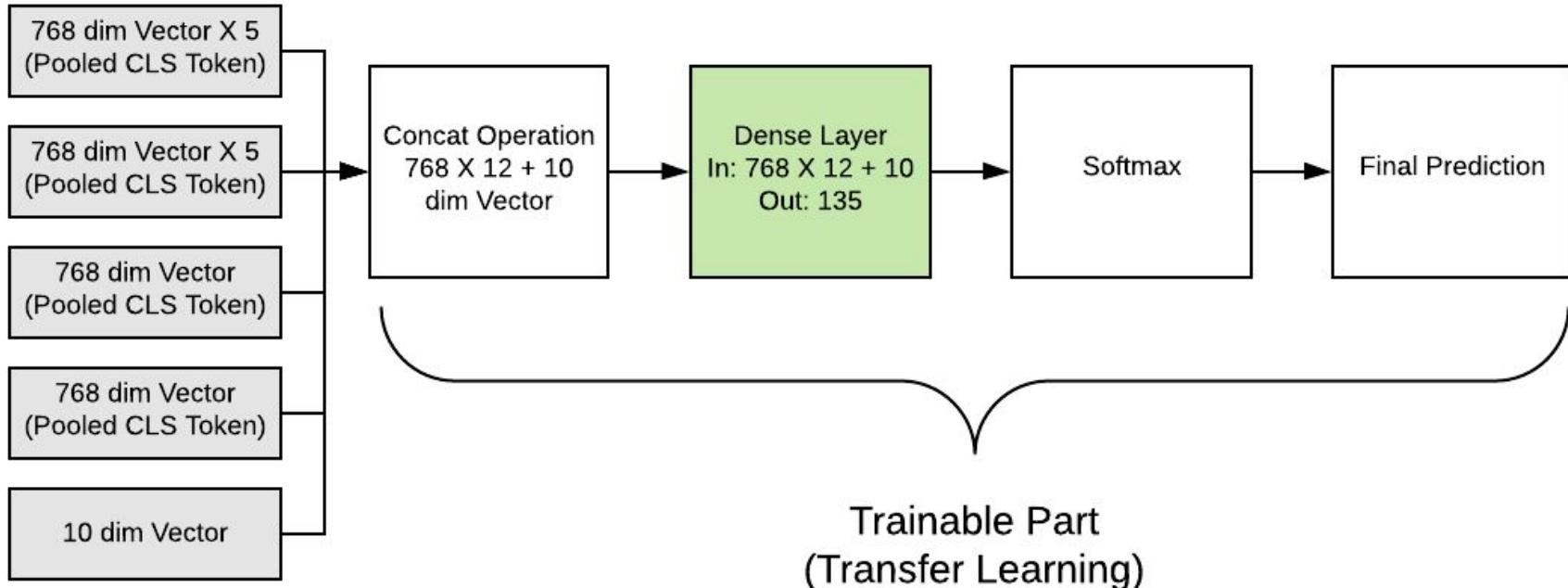
Implementation (moderate complexity)

- Use first 5 links per domain
- Free text:
 - Page Text
 - Page Title
 - Website Domain
- Metadata: Domain extension
- Other: Business Name

Model - Design



Model - Design



Model - Train

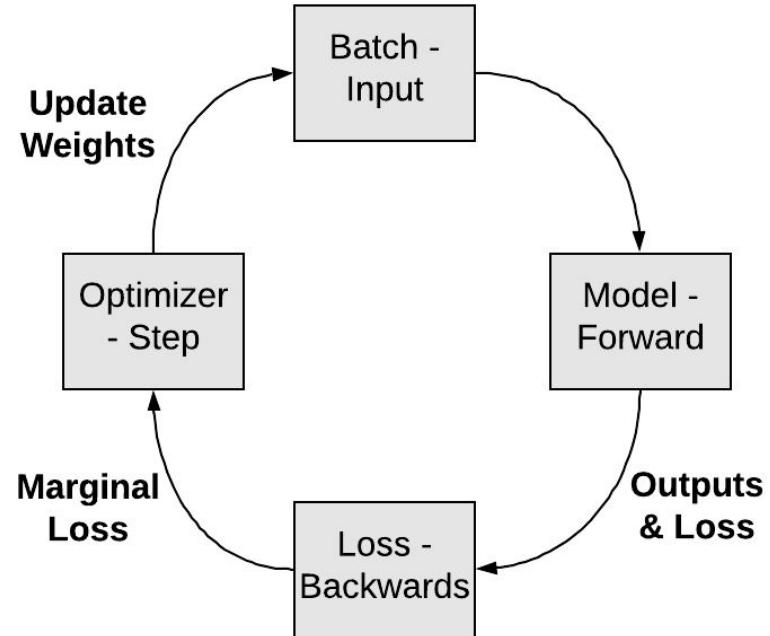
Tool:

- PyTorch
<https://github.com/pytorch/pytorch>
- Basically NumPy + GPUs
- Autograd
- Relatively easy to debug

Model - Train

NN / PyTorch Train Cycle:

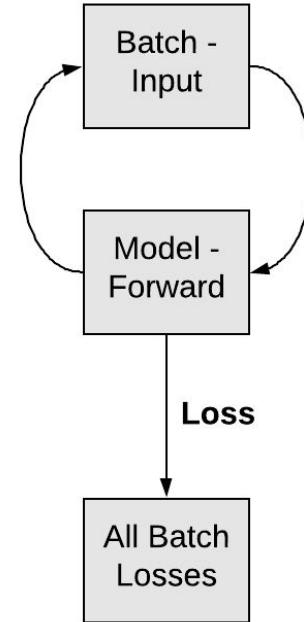
1. Select input batch (e.g. size 16)
2. Calc current outputs (forward)
3. Calc loss: Outputs VS labels
4. Calc marginal loss (backwards)
5. Update weights (step)
6. Again! Until batches / epoch done.



Model - Train

NN / PyTorch Validation Cycle:

1. Select input batch (e.g. size 16)
2. Calc current outputs (forward)
3. Calc loss: Outputs VS labels
4. Collect loss for all batches
5. Again! Until batches / epoch done
6. Report avg loss



Model - Train

Organizing the dataset:

- Convert all text elements into numeric vectors
- Convert everything to tensors
- Put everything together such that:
 - One “item” has all inputs per obs
 - Can be loaded in batches (GPU memory limits)

```
12 def make_dataset(input_df, bert_tokenizer):
13     """transform crawl data input into model readable tensor dataset..."""
14     # split input data into separate pages
15     num_pages = 5
16     data_pages = split_data_into_pages(input_df.copy(), num_pages)
17
18     # generate dummies for top 10 url extensions
19     top_extensions = ['biz', 'co', 'com', 'edu', 'gov',
20                         'info', 'net', 'org', 'site', 'us']
21     url_ext_mat = gen_url_extension_dummies(data_pages[0],
22                                              top_extensions)
23
24     # extract domain from url
25     data_pages[0]['url_domain_text'] = data_pages[0]['url'].apply(
26         extract_domain_from_url)
```



**DataFrame per
web page data**

```
34 # tokenize & truncate all input fields  
35 max_len = 256 ← Max num of tokens  
36  
37 page_text_tokenized = tokenize_and_truncate_sequence(  
38     data_pages,  
39     'response_text', ← Field to tokenize  
40     max_len,  
41     bert_tokenizer)  
42  
43 page_text_masks = gen_sequence_mask(page_text_tokenized)
```

**1 if token, 0
otherwise**

```
64  
65  
66  
67  
68  
69  
70  
71
```

```
# combine into tensor dataset  
data_set = combine_inputs_into_tensor_dataset(  
    page_text_tokenized, page_text_masks, page_titles_tokenized,  
    page_title_masks, domain_tokenized, domain_masks, biz_name_tokenized,  
    biz_name_masks, url_ext_mat)  
  
urls = data_pages[0]['url'].tolist()  
return data_set, urls
```



Big collection of tensors, ordered!

```
1 import torch
2 from torch import nn
3 from torch.utils.data import DataLoader, RandomSampler, SequentialSampler
4 from tqdm import trange
5 import numpy as np
6
7 from pytorch_transformers import BertTokenizer, AdamW
8
9 from industry_sagemaker_service.bert_custom_classifier import CustomBertClassifier
10 from industry_sagemaker_service.utils import make_dataset
```

**Loader for
BERT Tokenizer**



**Import the
make_dataset
method**

```
12 bert_model_name = 'bert-base-cased'  
13 bert_tokenizer = BertTokenizer.from_pretrained(bert_model_name,  
14                               do_lower_case=False)
```

**Load pretrained
BERT Tokenizer**

```
16 train_raw_data_df = '' # train crawl data + business names  
17 train_data = make_dataset(train_raw_data_df, bert_tokenizer)
```

**Make Train and
Test datasets**

```
19 test_raw_data_df = '' # test crawl data + business names  
20 test_data = make_dataset(test_raw_data_df, bert_tokenizer)
```

```
22 batch_size = 16
```

```
24 train_sampler = RandomSampler(train_data)  
25 train_dataloader = DataLoader(train_data, sampler=train_sampler,  
26                               batch_size=batch_size)
```

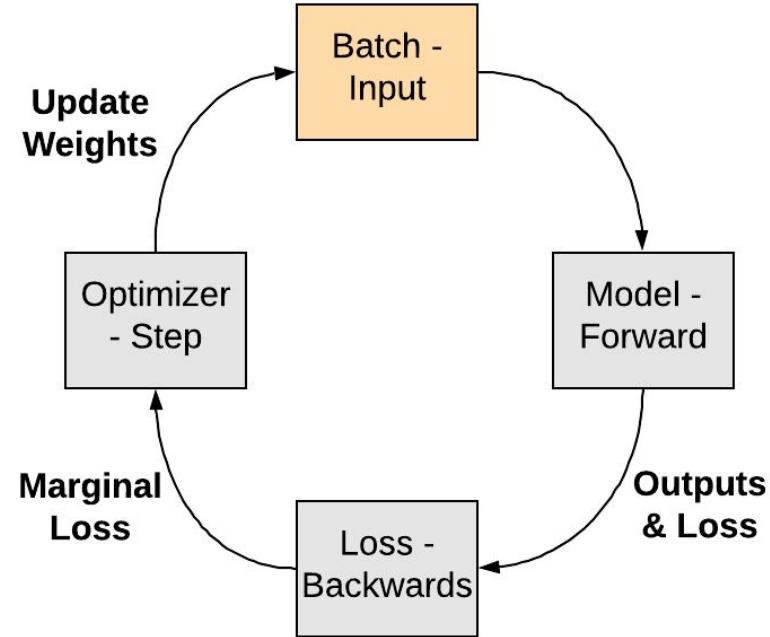
**Make Train
and Test
dataloaders**

```
28 test_sampler = SequentialSampler(test_data)  
29 test_dataloader = DataLoader(test_data, sampler=test_sampler,  
30                               batch_size=batch_size)
```

Model - Train

Batch - Input:

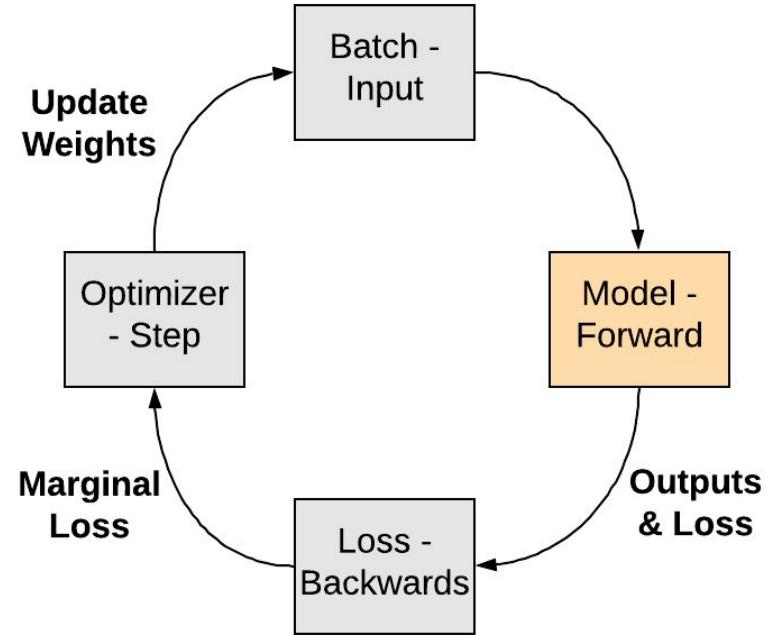
- Create Train Dataloader
- Extract Batch



Model - Train

Forward method:

- Define the graph of the model:
 - Layers
 - Transformations
 - Combinations
- Is invoked with inputs, calculates outputs
- Using labels and a loss metric, calculates loss



```
1 from pytorch_transformers import BertModel, BertPreTrainedModel
2 from torch.nn import CrossEntropyLoss, MSELoss
3 from torch import nn
4 import torch
5
6
7 class CustomBertClassifier(BertPreTrainedModel):
8     def __init__(self, config):
9         super(CustomBertClassifier, self).__init__(config)
10        self.num_labels = config.num_labels
11        self.bert = BertModel(config)
12        self.dropout = nn.Dropout(config.hidden_dropout_prob)
13        self.classifier = nn.Linear(config.hidden_size * 12 + 10,
14                                    self.config.num_labels)
15
16        self.init_weights()
```

Pre-Trained
BERT

Trainable
Linear Layer

```
18 ⬆  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36
```

```
def forward(self,  
           sentence_ids, sentence_masks,  
           title_ids, title_masks,  
           domain_ids, domain_masks,  
           biz_name_ids, biz_name_masks,  
           url_exts,  
           token_type_ids=None,  
           position_ids=None, head_mask=None,  
           labels=None):
```

```
# domain input  
domain_outputs = self.bert(  
    domain_ids,  
    attention_mask=domain_masks,  
    token_type_ids=token_type_ids,  
    position_ids=position_ids,  
    head_mask=head_mask)
```

```
domain_pooled_output = domain_outputs[1]
```

**All inputs:
Tokens + masks
& dummies**

**Domain Tokens
+ masks fed
into BERT**

**Domain
Embedding Vector
Obtained**

```
48 # page text inputs  
49 page_text_outputs = []  
50 for page_text_id, page_text_mask in \  
51     zip(page_text_ids, page_text_masks):  
52         page_text_output_tmp = self.bert(  
53             page_text_id,  
54             attention_mask=page_text_mask,  
55             token_type_ids=token_type_ids,  
56             position_ids=position_ids,  
57             head_mask=head_mask)  
58  
59         page_text_pooled_output_tmp = page_text_output_tmp[1]  
60         page_text_outputs += [page_text_pooled_output_tmp]
```

**Iterating over all
5 page texts**

**Page texts
Tokens + masks
fed into BERT**

**Page Text
embeddings
Obtained**

```
76  
77 # combine all outputs  
78 outputs_all = page_text_outputs + page_title_outputs + \  
79     [domain_pooled_output, biz_name_pooled_output, url_exts]  
80 pooled_output = torch.cat(outputs_all, dim=1)  
81 pooled_output = self.dropout(pooled_output)
```

```
82  
83 # apply classifier
```

```
logits = self.classifier(pooled_output)
```



**Big vector sent
through the
trainable Linear
layer**



**All vectors
concatenated into
a big vector**

```
85 # add hidden states and attention if they are here  
86 outputs = (logits,)  
87  
88 if labels is not None:  
89     if self.num_labels == 1:  
90         # We are doing regression  
91         loss_fct = MSELoss()  
92         loss = loss_fct(logits.view(-1), labels.view(-1))  
93     else:  
94         loss_fct = CrossEntropyLoss()  
95         loss = loss_fct(logits.view(-1, self.num_labels),  
96                             labels.view(-1))  
97  
98         outputs = (loss,) + outputs  
99 return outputs
```

Return forward result

**In Train mode: also return
loss**

**In Train mode:
Loss = Batch
output VS
Labels**

**Import and init
the custom
classifier**



```
33 from industry_sagemaker_service.bert_custom_classifier \
34     import CustomBertClassifier
35
36 custom_bert_mdl = CustomBertClassifier.from_pretrained(bert_model_name,
37                                         num_labels=135)
```

**Setup for
multiple GPUs**



```
39 if torch.cuda.device_count() > 1:
40     print("Let's use", torch.cuda.device_count(), "GPUs!")
41     custom_bert_mdl = nn.DataParallel(custom_bert_mdl)
42
43 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
44 custom_bert_mdl.to(device)
```

```
69 train_loss_set = []
70 epochs = 3
71
72 for _ in trange(epochs, desc="Epoch"):
73     custom_bert_mdl.train()
74
75     tr_loss, tr_accuracy = 0, 0
76     nb_tr_examples, nb_tr_steps = 0, 0
77
78     for step, batch in enumerate(train_dataloader):
79         # Add batch to GPU
80         batch = tuple(t.to(device) for t in batch)
81
82         # Unpack the inputs from our dataloader
83         page_text_ids1, page_text_ids2, page_text_ids3, \
84         page_text_ids4, page_text_ids5, \
85             page_text_masks1, page_text_masks2, page_text_masks3, \
86             page_text_masks4, page_text_masks5, \
87             page_title_ids1, page_title_ids2, page_title_ids3, \
88             page_title_ids4, page_title_ids5, \
89             page_title_masks1, page_title_masks2, page_title_masks3, \
90             page_title_masks4, page_title_masks5, \
91             domain_ids, domain_masks, biz_name_ids, biz_name_masks, \
92             url_exts, batch_labels = batch
```

Allow model to update its weights

Iterate over all batches in train dataloader

Unpack all inputs from batch

Loss & Outputs calculated

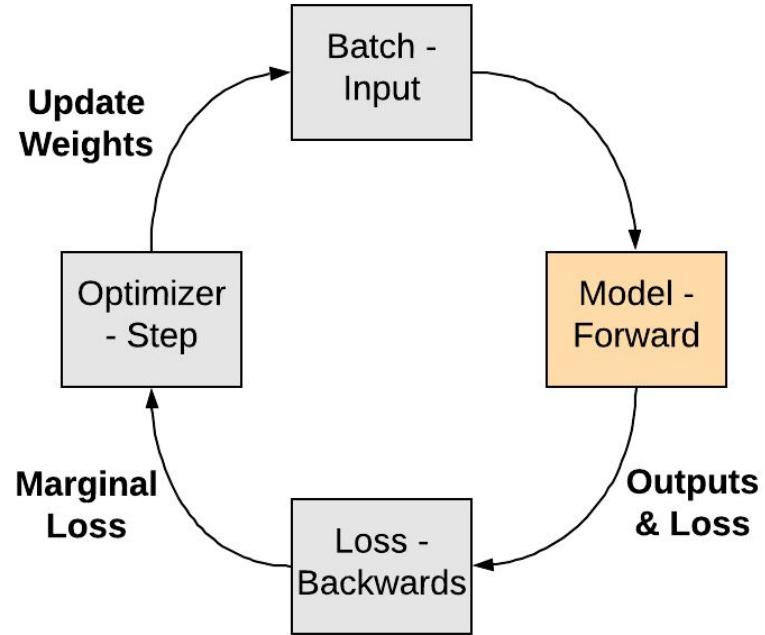
All inputs sent
to forward
method

```
96  
97  
98  
99  
100 # Clear out the gradients (by default they accumulate)  
optimizer.zero_grad()  
101  
102  
103  
104  
105  
106  
107  
108  
109  
110  
111  
112  
113  
114  
115 # Forward pass  
loss, logits = custom_bert_mdl(  
    sentence_ids=(page_text_ids1, page_text_ids2, page_text_ids3,  
                  page_text_ids4, page_text_ids5),  
    sentence_masks=(page_text_masks1, page_text_masks2, page_text_masks3,  
                  page_text_masks4, page_text_masks5),  
    title_ids=(page_title_ids1, page_title_ids2, page_title_ids3,  
              page_title_ids4, page_title_ids5),  
    title_masks=(page_title_masks1, page_title_masks2, page_title_masks3,  
              page_title_masks4, page_title_masks5),  
    domain_ids=domain_ids,  
    domain_masks=domain_masks,  
    biz_name_ids=biz_name_ids,  
    biz_name_masks=biz_name_masks,  
    url_exts=url_exts,  
    token_type_ids=None,  
    labels=batch_labels)
```

Model - Train

Forward method:

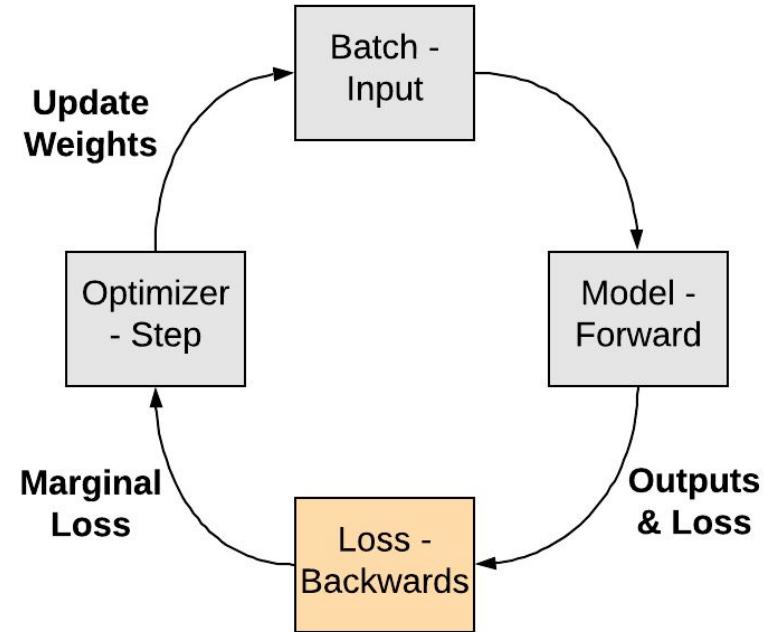
- Define the graph of the model:
 - Layers
 - Transformations
 - Combinations
- Is invoked with inputs, calculates outputs
- Using labels and a loss metric, calculates loss



Model - Train

Loss metric:

- Compares outputs with labels, calculates actual loss
- “backward” method computes the marginal loss w.r.t all the trainable parameters (i.e the loss gradient). Updates “grad” attribute of tensors.
- E.g: Cross Entropy Loss for multi-class problems



```
96      # Clear out the gradients (by default they accumulate)
97      optimizer.zero_grad()
98
99      # Forward pass
100     loss, logits = custom_bert_mdl(
101         sentence_ids=(...),
103         sentence_masks=(...),
105         title_ids=(...),
107         title_masks=(...),
109         domain_ids=domain_ids,
110         domain_masks=domain_masks,
111         biz_name_ids=biz_name_ids,
112         biz_name_masks=biz_name_masks,
113         url_exts=url_exts,
114         token_type_ids=None,
115         labels=batch_labels)
```

**Aggregate
loss from
multiple GPUs**



```
118     train_loss_set.append(loss.mean().item())
```

```
119     loss.mean().backward()
120
```

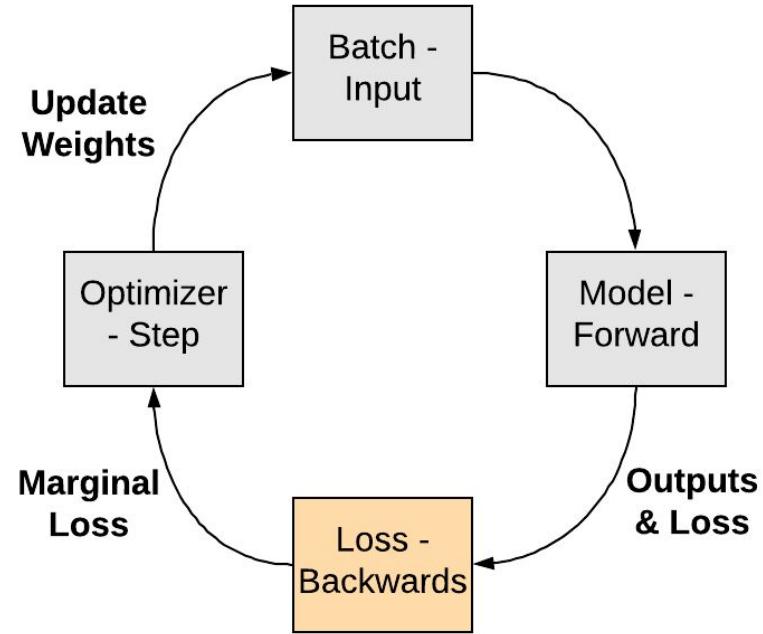
**Calculate
Marginal Loss**



Model - Train

Loss metric:

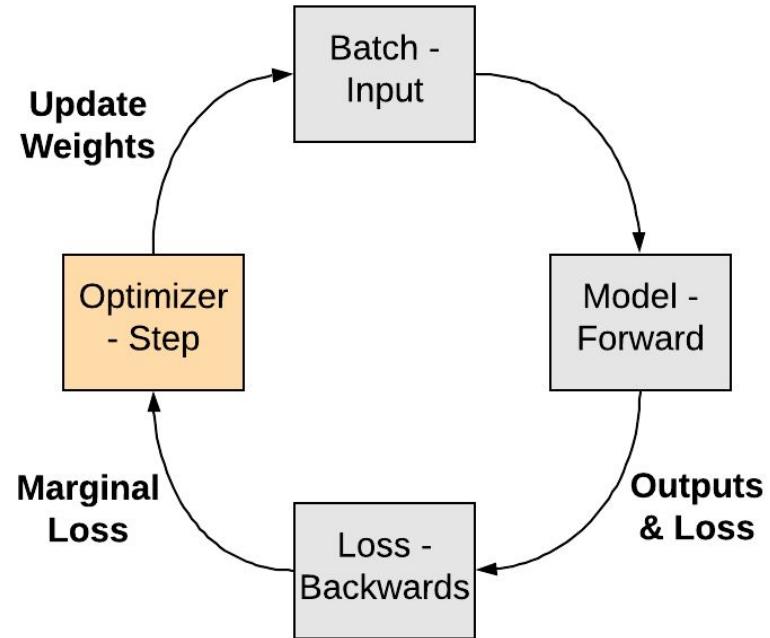
- Compares outputs with labels, calculates actual loss
- “backward” method computes the marginal loss w.r.t all the trainable parameters (i.e the loss gradient). Updates “grad” attribute of tensors.
- E.g: Cross Entropy Loss for multi-class problems



Model - Train

Optimizer:

- Specifies which gradient-based optimization method to use for gradient descent
- “step” method updates all relevant weights (i.e. tensor values) according to the marginal loss (stored in the grad attribute of the tensors)
- E.g: AdamW Optimizer



```
47 param_optimizer = list(custom_bert_mdl.named_parameters())
48 no_decay = ['bias', 'gamma', 'beta']
49 optimizer_grouped_parameters = [
50     {'params': [p for n, p in param_optimizer if
51                 not any(nd in n for nd in no_decay)],
52      'weight_decay_rate': 0.01},
53     {'params': [p for n, p in param_optimizer if
54                 any(nd in n for nd in no_decay)],
55      'weight_decay_rate': 0.0}
56 ]
57
58 optimizer = AdamW(optimizer_grouped_parameters, lr=2e-5)
```

```
96 # Clear out the gradients (by default they accumulate)
97 optimizer.zero_grad()
98
99 # Forward pass
100 loss, logits = custom_bert_mdl(
101     sentence_ids=(...),
102     sentence_masks=(...),
103     title_ids=(...),
104     title_masks=(...),
105     domain_ids=domain_ids,
106     domain_masks=domain_masks,
107     biz_name_ids=biz_name_ids,
108     biz_name_masks=biz_name_masks,
109     url_exts=url_exts,
110     token_type_ids=None,
111     labels=batch_labels)
112
113
114
115
116
117
118 train_loss_set.append(loss.mean().item())
119
120 loss.mean().backward()
121
122 optimizer.step()
```

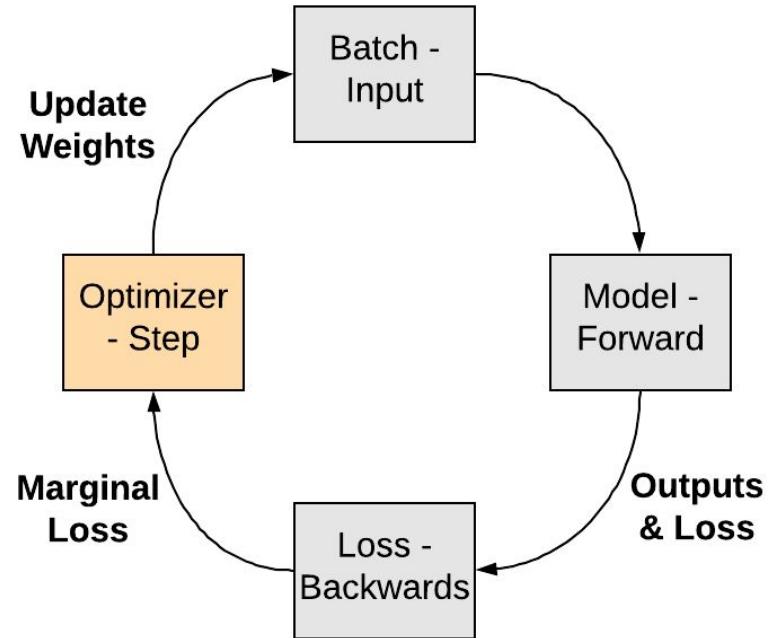
**Clear gradients
from previous
batches**

**Update weights
based on
marginal loss**

Model - Train

Optimizer:

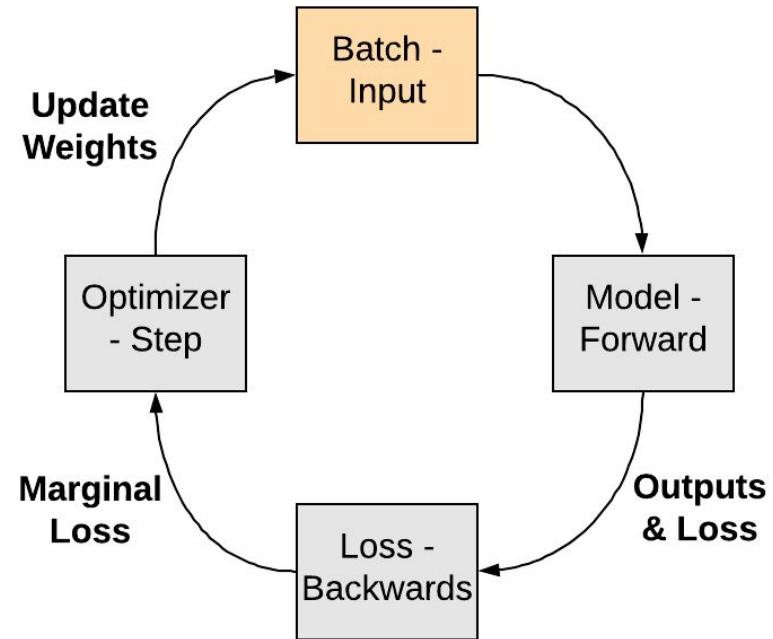
- Specifies which gradient-based optimization method to use for gradient descent
- “step” method updates all relevant weights (i.e. tensor values) according to the marginal loss (stored in the grad attribute of the tensors)
- E.g: AdamW Optimizer



Model - Train

Batch - Input:

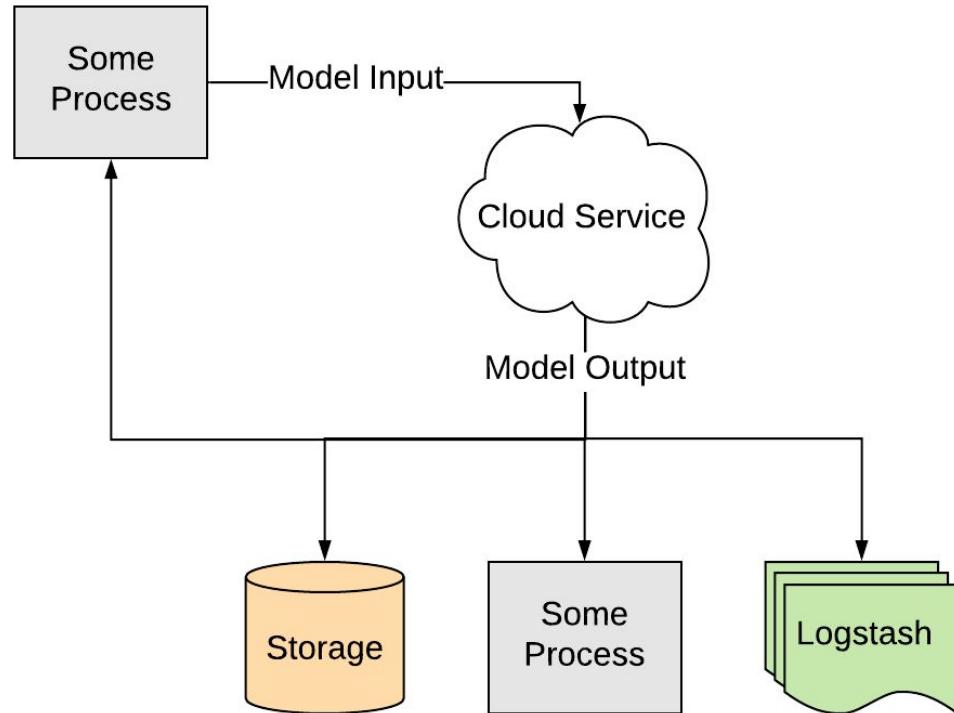
- Next Batch



Model - Deploy

Goal:

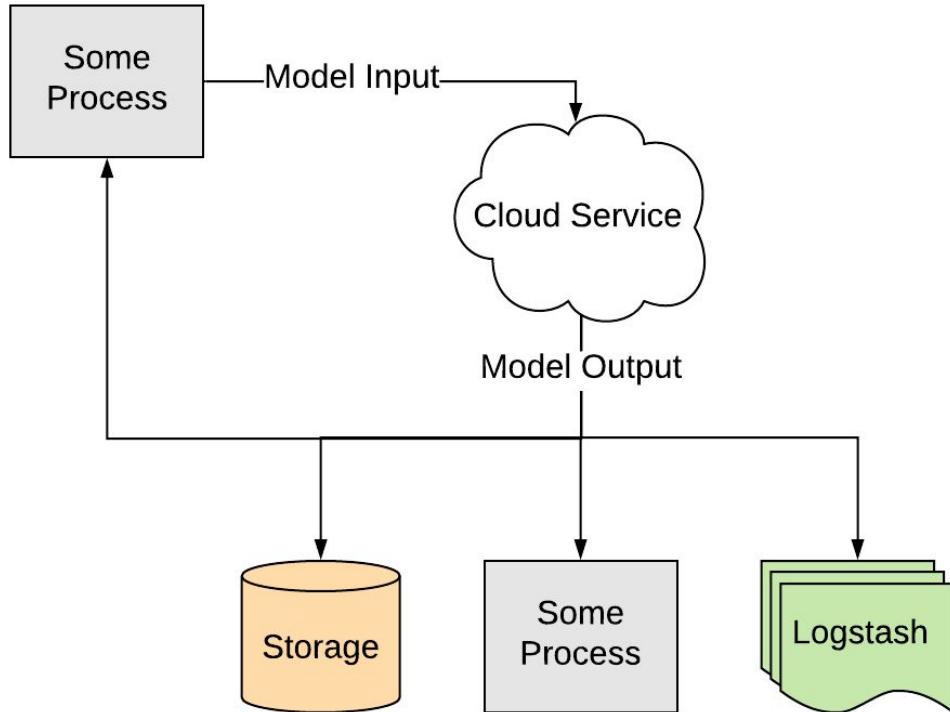
- Model hosted in the cloud
- Always up (persistence)
- Communicates via API
- Can scale



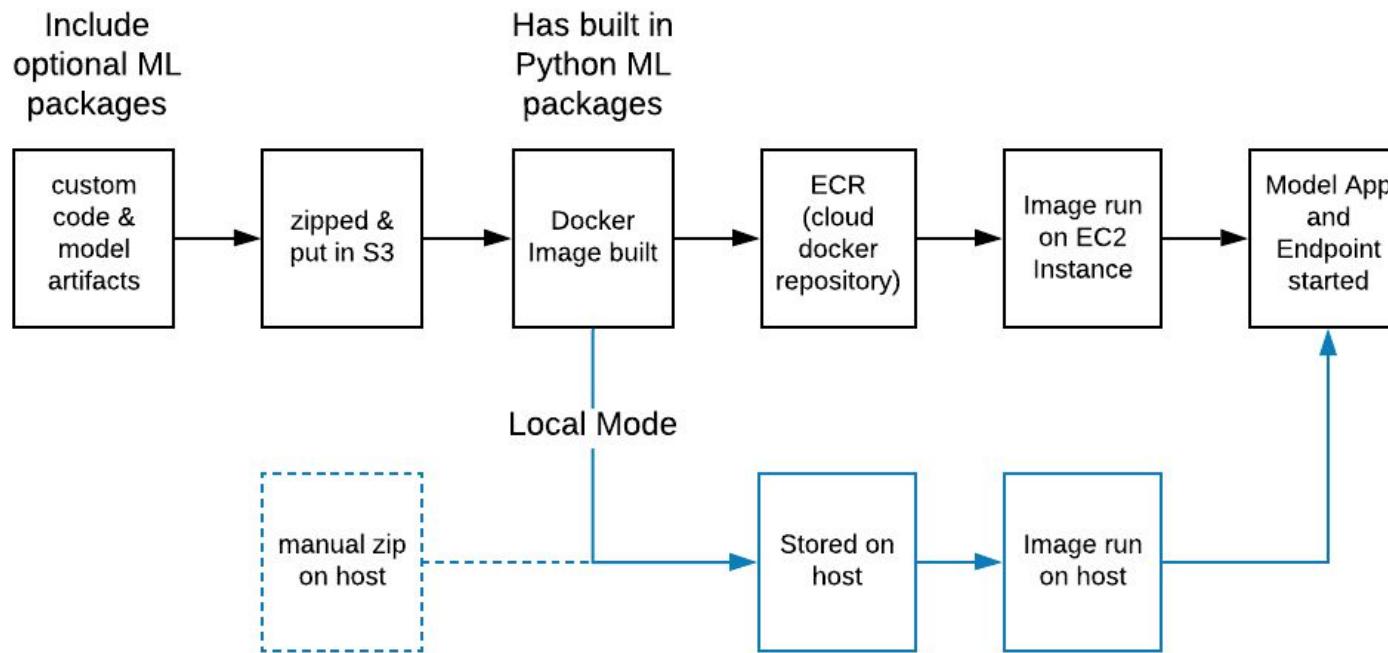
Model - Deploy

Tool:

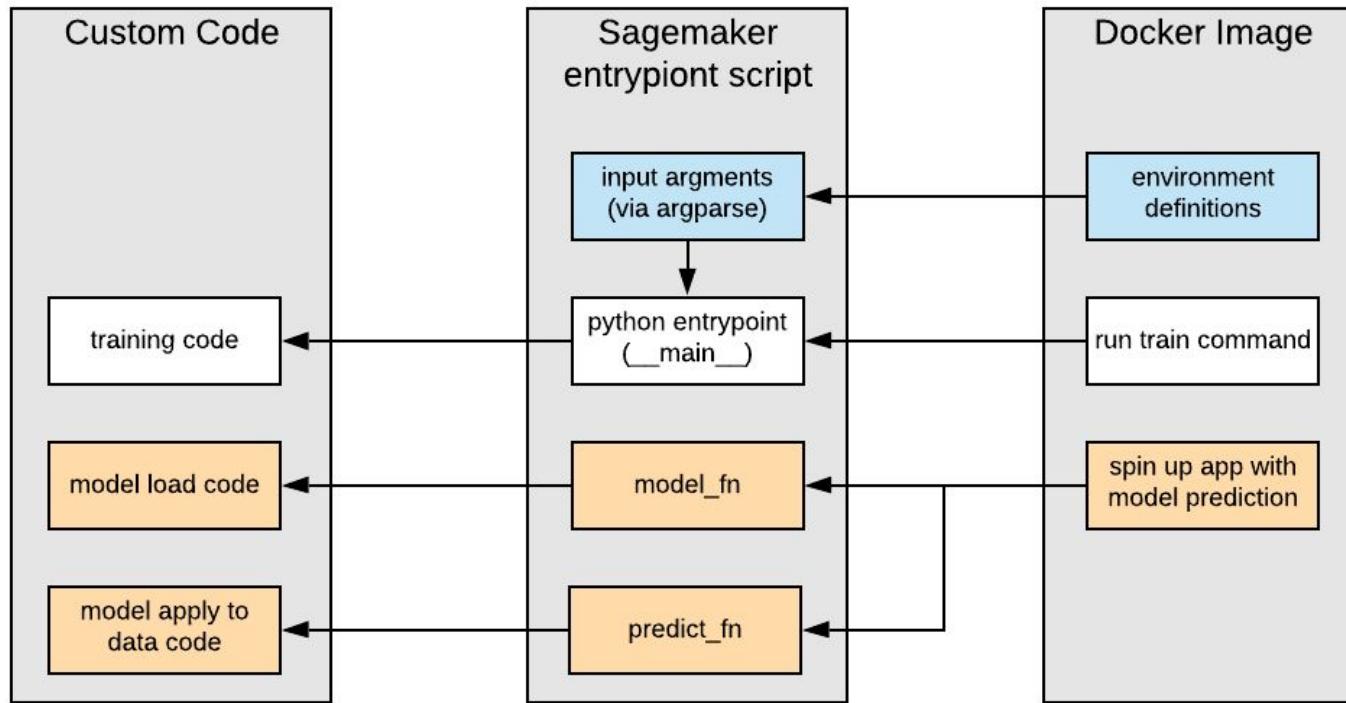
- Amazon SageMaker
<https://aws.amazon.com/sagemaker/>
- All the pros of docker deployment without writing any docker code
- Supports four modes: Tag, Explore, Train and Deploy
- Auto Scaling



Model - Deploy (SageMaker)



Model - Deploy (SageMaker)



```
14 def model_fn(model_dir):  
15     """  
16         initialize resources during service startup  
17     """  
18     # load model  
19     from industry_sagemaker_service.bert_custom_classifier \  
20         import CustomBertClassifier  
21     bert_mdl = CustomBertClassifier.from_pretrained(  
22         os.path.join(model_dir, 'data_files', 'model_e'))  
23  
24     # set up for GPU if exists  
25     device = torch.device("cuda" if torch.cuda.is_available() else "cpu")  
26     if device.type == 'cuda':  
27         bert_mdl.cuda()  
28  
29     # load tokenizer  
30     bert_tokenizer = BertTokenizer.from_pretrained(  
31         'bert-base-cased', do_lower_case=False)  
32  
33     # load industry map, int --> industry name  
34     industry_map_df = pd.read_csv(os.path.join(  
35         model_dir, 'data_files', 'industry_map.csv'))  
36     industry_map_dict = industry_map_df.\\  
37         set_index('industry_int').to_dict()['industry_tag']  
38  
39     return bert_mdl, bert_tokenizer, industry_map_dict
```

Load trained model

Load pertained tokenizer

Return all artifacts

```
73 def predict_fn(crawl_df, model_fn_input): ← Get preloaded model  
74     """  
75         apply model to input, return inference  
76     """  
77     bert_mdl, bert_tokenizer, industry_map_dict = model_fn_input  
78  
79     from industry_sagemaker.service.utils import make_dataset, predict_industry  
80     data_set, urls = make_dataset(crawl_df, bert_tokenizer)  
81     data_sampler = SequentialSampler(data_set)  
82     data_loader = DataLoader(data_set, sampler=data_sampler, batch_size=8) ← Turn input into  
83  
84     device = torch.device("cuda" if torch.cuda.is_available() else "cpu")  
85  
86     preds = []  
87     probas = []  
88     for step, batch in enumerate(data_loader):... ← Inference loop by batch  
89  
90     preds_as_str = [[industry_map_dict[y] for y in x] for x in preds]  
91     probas_rounded = [[np.round(y, 4) for y in x] for x in probas]  
92     output = [{k: v for (k, v) in zip(x, y)} for (x, y) in  
93             zip(preds_as_str, probas_rounded)]  
94     output_sorted = [sorted(d.items(), key=operator.itemgetter(1),  
95                           reverse=True) for d in output] ← Sort preds in  
96  
97     assert len(urls) == len(output_sorted)  
98     output_with_urls = {x: y for (x, y) in zip(urls, output_sorted)}  
99  
100    return json.dumps(output_with_urls)
```

```
1 import sagemaker  
2 from boto3 import Session as BotoSession  
3  
4 # create sagemaker session  
5 boto_session = BotoSession(profile_name='default',  
6                             region_name='us-east-1')  
7 sagemaker_session = sagemaker.Session(boto_session=boto_session)  
8 sagemaker_role = 'ido-sagemaker-test'  
9  
10 from sagemaker.pytorch import PyTorchModel  
11  
12 industry_cloud_model = PyTorchModel(  
13  
14     entry_point='sagemaker_entry_point.py',  
15  
16     model_data='s3://ido-sagemaker-test/ind_mdl_data.tar.gz',  
17  
18     framework_version='1.1.0',  
19     py_version='py3',  
20  
21     role=sagemaker_role,  
22     source_dir='.')
```

Authenticate with SageMaker service

Specify SageMaker PyTorch Docker Container

Location of entrypoint script

Location of model artifacts in S3

Specify PyTorch and Python versions

Call the deploy method on the model

```
24 industry_cloud_predictor = industry_cloud_model.deploy(  
25     initial_instance_count=1,  
26     instance_type='ml.p2.xlarge')
```

Specify EC2 Instance (GPU)

```
27  
28  
29 from sagemaker.predictor import json_serializer  
30  
31 # define input format as JSON  
32 industry_cloud_predictor.serializer = json_serializer  
33 industry_cloud_predictor.content_type = 'application/json'  
34  
35 input_json = '' # load some input data  
36  
37 preds = industry_cloud_predictor.predict(input_json)
```

Call the predict method on the predictor

Endpoints

Update endpoint

Actions ▾

Create endpoint

Search endpoints

< 1 > ⚙

Our Endpoint

↓

Name	ARN	Creation time	Status	Last updated
sagemaker-pytorch-2019-10-30-03-50-29-502	[REDACTED]	Oct 30, 2019 03:50 UTC	✓ InService	Oct 30, 2019 04:00 UTC
sagemaker-scikit-learn-2019-10-08-21-54-08-741	[REDACTED]	Oct 08, 2019 21:54 UTC	✓ InService	Oct 25, 2019 19:34 UTC

Conclusion - Summary

What we did:

- Built a 250K sized sample of websites
 - Search strategies
 - Balancing strategies
- Tagged that sample
 - Crowdsourcing considerations & strategies
- Crawled all of the websites in it
 - Scrapy for small scale
 - ScrapyD for large scale

Conclusion - Summary

What we did:

- Built a 250K sized sample of websites
 - Search strategies
 - Balancing strategies
- Tagged that sample
 - Crowdsourcing considerations & strategies
- Crawled all of the websites in it
 - Scrapy for small scale
 - ScrapyD for large scale
- Built an NLP industry classification model
 - Transfer learning using BERT
 - Dealt with 512 length limit
- Trained that model with multiple GPUs
 - Pretrained BERT models with PyTorch
- Deployed that model using SageMaker
 - Docker deployment, no docker code

Conclusion - Key Takeaways

- Crowdsourcing has to be done carefully to produce real value
 - Vendor trials
 - Enforce QA
- You can build a very efficient crawler using Scrapy / ScrapyD
- You can build a very powerful NLP model using pretrained BERT
 - Use available wrappers (e.g HuggingFace for PyTorch)
- You can train / deploy models on very powerful machines (e.g. 16 GPUs) using SageMaker
 - No docker knowledge required
- Hope for the best



Questions? + Thanks!

ido.shlomo@bluevine.com