

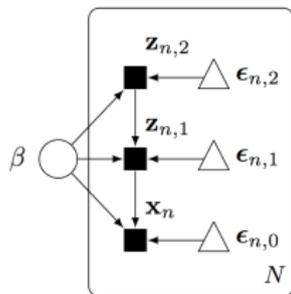
What Might Deep Learners Learn From Probabilistic Programming?

Dustin Tran
Google Brain

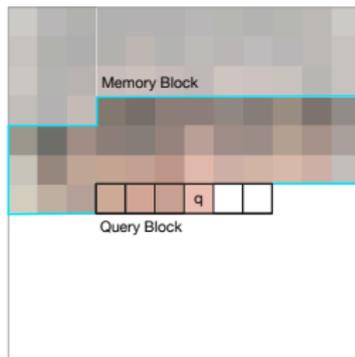


Interested in research for scientific applications?

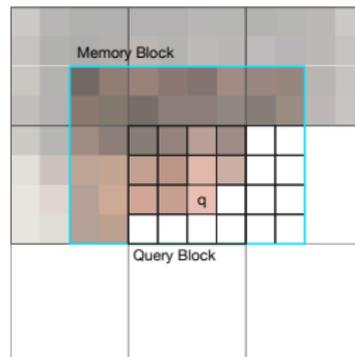
That's not this talk.



Local 1D Attention



Local 2D Attention



Generation & compression of 10M colored 32x32 images

23 Million Parameters

title: "Pescadores Weekly" length: 2000

Pescadores Weekly Inc (organized in 1978 as the Pescadores Weekly) is the official Weekly organizing body of science. It was established in 1978 atop the city of the center of the City of Gravity. In 1978 it became the United States National Museum of Science and Design. The growing Berklee Plant represents the physical community of typical worldwide plants.

Pescadores is named after the lord of the University of the Gravity of Drama, an interest he established in 1989 of the West Bird Plant. In the fall of 1978 the Faxophonius Institute, with the Today Natural History Foundation, published by the University of Georgia, and Anubis (Yale University).

==History==

Pescadores provided classic books on science and other philosophical sciences such as geoscience that accomplished science-fiction in Europe during 1976–1977, being the first university campus of the University of New Zealand. ...

340 Million Parameters

title: "Pescadores Weekly" length: 2000

The "Pescadores Weekly" is the only daily newspaper in the Fort Hood metropolitan area, although the entire population of Fort Hood is majority white in population. Due to its not-for-profit focus as a newspaper, Carehead reporters and Hometown have expressed the need for a feature to focus on Fort Hood and its inner-city population.

According to the Pescadores Weekly Editors' annual survey, the "Pescadores Weekly" circulation was the largest in Texas, surpassing the "Evening Star Sans Monthly" made in 1992. By contrast, Terry Truehead's "The Best of Fort Hood" daily was the fourth-fastest in Texas (behind rival "AfterEllena" and "Expensive Ellena", topping the chart thirteen years in a row) and expanded by two thirds during the 2002-2003 Edition.

==Political affiliation==

The newspaper endorsed John McCain's majoritarian, Jack Abramoff, at the 2004 GOP convention. ...

Scaling up fundamental language models

[Liu+ 2018; Shazeer+ 2018]

Inference in a probabilistic program

```
(trace, weight) = query(program, args, observations)
```

\mathbf{x}

ξ

\mathcal{P}

α

\mathbf{y}

The Myth of Probabilistic Programming

**Programming is infeasible if a core operation
in the language is NP-hard.**

For high-dimensional problems + modern probabilistic models, we haven't solved automated inference.

<> Code

🔔 Issues 117

🔀 Pull requests 23

📊 Insights

A library for probabilistic modeling, inference, and criticism. Deep generative models, variational inference. Runs on TensorFlow. <http://edwardlib.org>

bayesian-methods

deep-learning

machine-learning

data-science

tensorflow

neural-networks

statistics

probabilistic-programming

📄 1,761 commits

🌿 19 branches

📦 27 releases

👤 66 contributors

Branch: master ▾

New pull request

Find file

Clone or download ▾



christopherlovell committed with dustinvtran fixed invgamma_normal_mh example (#793) ⋮

Latest commit 081ea53 23 days ago



docker Use Observations and remove explicit storage of data files (#751)

3 months ago



docs Revise docs to enable spaces in filepaths; update travis with tf==1.4...

26 days ago



Sign Up

Log In



all categories ▾

Latest

Top

Topic	Category	Users	Replies	Views	Activity
Iterative estimators ("bayes filters") in Edward?			5	21	7h
Tutorial for multiple variational methods using Poisson regression?			2	20	1d



blei-lab/edward

A library for probabilistic modeling, inference, and criticism. <http://edwardlib.org>

Faez Shakil @faezs

Jan 23 02:47

Hi @dustinvtran, thanks for edward, the library and surrounding literature have been immense fun to get into. Would you be able to tell me whether it'd be relatively painless to get the inference compute graphs from Ed as native tensorflow graphdefs and use them on mobile platforms? Or would I have to port a bunch of custom ops from edward into my own tensorflow build for use on mobile?

PEOPLE REPO INFO



Edward

Failure Modes

- **Inference is monolithic.** The average workflow requires understanding a new ecosystem, closed under its own compositions.
- **Can't it go faster?** Edward was not designed with TPUs and multiple machines in mind.

Some Iteration of Edward

Random Variables Are All You Need

Edward2 reifies any computable probability distribution as a Python function. Inputs to the program represent values the distribution conditions on.

```
def model():
    p = Beta(1., 1., name="p")
    x = Bernoulli(probs=p, sample_shape=50, name="x")
    return x

-

import neural_net_negative, neural_net_positive

def variational(x):
    eps = Normal(0., 1., sample_shape=2, name="eps")
    if eps[0] > 0:
        return neural_net_positive(eps[1], x)
    else:
        return neural_net_negative(eps[1], x)
```

Tracing

A tracer from AD wraps the language's primitive operations. The tracer intercepts control just before those operations are executed.

Edward2 applies tracing in order to perform user-programmable manipulations.

```
INTERCEPTOR_STACK = [lambda f, *args, **kwargs: f(*args, **kwargs)]

@contextmanager
def interception(interceptor):
    INTERCEPTOR_STACK.append(interceptor)
    yield
    INTERCEPTOR_STACK.pop()

def interceptable(func):
    def func_wrapped(*args, **kwargs):
        INTERCEPTOR_STACK[-1](func, *args, **kwargs)
    return func_wrapped
```

Example: Latent Dirichlet Allocation

```
1 from __future__ import absolute_import
2 from __future__ import division
3 from __future__ import print_function
4
5 import functools
6 import os
7
8 from absl import flags
9 import numpy as np
10 import scipy.sparse
11 from six.moves import cPickle as pickle
12 from six.moves import urllib
13 import tensorflow as tf
14
15 from tensorflow_probability import edward2 as ed
16
17
18 flags.DEFINE_float(
19     "learning_rate", default=3e-4, help="Learning rate.")
20 flags.DEFINE_integer(
21     "max_steps", default=180000, help="Number of training steps to run.")
22 flags.DEFINE_integer(
23     "num_topics",
24     default=50,
25     help="The number of topics.")
26 flags.DEFINE_list(
27     "layer_sizes",
28     default=["300", "300", "300"],
29     help="Comma-separated list denoting hidden units per layer in the encoder.")
30 flags.DEFINE_string(
31     "activation",
32     default="relu",
33     help="Activation function for all hidden layers.")
34 flags.DEFINE_integer(
35     "batch_size",
36     default=32,
37     help="Batch size.")
38 flags.DEFINE_float(
39     "prior_initial_value", default=0.7, help="The initial value for prior.")
40 flags.DEFINE_integer(
41     "prior_burn_in_steps",
42     default=120000,
43     help="The number of training steps with fixed prior.")
44 flags.DEFINE_string(
45     "data_dir",
46     default=os.path.join(os.getenv("TEST_TMPDIR"), "/tmp"), "lda/data",
47     help="Directory where data is stored (if using real data).")
48 flags.DEFINE_string(
49     "model_dir",
50     default=os.path.join(os.getenv("TEST_TMPDIR"), "/tmp"), "lda/",
51     help="Directory to put the model's fit.")
52 flags.DEFINE_integer(
53     "viz_steps", default=10000, help="Frequency at which save visualizations.")
54 flags.DEFINE_bool("fake_data", default=False, help="If true, uses fake data.")
55 flags.DEFINE_bool(
56     "delete_existing",
57     default=False,
58     help="If true, deletes existing directory.")
59
60 FLAGS = flags.FLAGS
61
62
63 def clip_dirichlet_parameters(x):
64     """Clips the Dirichlet parameters to the numerically stable KL region."""
65     return tf.clip_by_value(x, 1e-3, 1e3)
66
67
68 def latent_dirichlet_allocation(concentration, topics_words):
69     topics = ed.Dirichlet(concentration=concentration, name="topics")
70     word_probs = tf.matmul(topics, topics_words)
71     # The observations are bags of words and therefore not one-hot. However,
72     # log_prob of OneHotCategorical computes the probability correctly in
73     # this case.
74     bag_of_words = ed.OneHotCategorical(probs=word_probs, name="bag_of_words")
75     return bag_of_words
76
77
78 def make_lda_variational(activation, num_topics, layer_sizes):
79     encoder_net = tf.keras.Sequential()
80     for num_hidden_units in layer_sizes:
81         encoder_net.add(tf.keras.layers.Dense(
82             num_hidden_units, activation=activation,
83             kernel_initializer=tf.glorot_normal_initializer()))
84     encoder_net.add(tf.keras.layers.Dense(
85         num_topics, activation=tf.nn.softplus,
86         kernel_initializer=tf.glorot_normal_initializer()))
87
88 def lda_variational(bag_of_words):
89     concentration = _clip_dirichlet_parameters(encoder_net(bag_of_words))
90     return ed.Dirichlet(concentration=concentration, name="topics_posterior")
91
92
93 def model_fn(features, labels, mode, params, config):
94     del labels, config
95
96     # Set up the model's learnable parameters.
97     logit_concentration = tf.get_variable(
98         "logit_concentration",
99         shape=[1, params["num_topics"]],
100         initializer=tf.constant_initializer(
101             _softplus_inverse(params["prior_initial_value"]]))
102     concentration = _clip_dirichlet_parameters(
103         tf.nn.softplus(logit_concentration))
104
105     num_words = features.shape[1]
106     topics_words_logits = tf.get_variable(
107         "topics_words_logits",
108         shape=[params["num_topics"], num_words],
109         initializer=tf.glorot_normal_initializer())
110     topics_words = tf.nn.softmax(topics_words_logits, axis=-1)
111
112     # Compute expected log-likelihood. First, sample from the variational
113     # distribution; second, compute the log-likelihood given the sample.
114     lda_variational = make_lda_variational(
115         params["activation"],
116         params["num_topics"],
117         params["layer_sizes"])
118     with ed.tape() as variational_tape:
119         _ = lda_variational(features)
120
121     with ed.tape() as model_tape:
122         with ed.interception(
123             make_value_setter(topics=variational_tape["topics_posterior"])):
124             posterior_predictive = latent_dirichlet_allocation(concentration,
125                 topics_words)
126
127     log_likelihood = posterior_predictive.distribution.log_prob(features)
128     tf.summary.scalar("log_likelihood", tf.reduce_mean(log_likelihood))
129
130     # Compute the KL-divergence between two Dirichlets analytically.
131     # The sampled KL does not work well for "sparse" distributions
```

Example: Latent Dirichlet Allocation

```
144 # Ensure that the KL is non-negative (up to a very small slack).
145 # Negative KL can happen due to numerical instability.
PROBPROG 2018 Schedule | PROBPROG 2018 [t.greater(kl, -1e-3, message="kl!")]:
148
149 elbo = log_likelihood - kl
150 avg_elbo = tf.reduce_mean(elbo)
151 tf.summary.scalar("elbo", avg_elbo)
152 loss = -avg_elbo
153
154 # Perform variational inference by minimizing the -ELBO.
155 global_step = tf.train.get_or_create_global_step()
156 optimizer = tf.train.AdamOptimizer(params["learning_rate"])
157
158 # This implements the "burn-in" for prior parameters (see Appendix D of [2]
159 # For the first prior_burn_in_steps steps they are fixed, and then trained
160 # jointly with the other parameters.
161 grads_and_vars = optimizer.compute_gradients(loss)
162 grads_and_vars_except_prior = [
163     x for x in grads_and_vars if x[1] != logit_concentration]
164
165 def train_op_except_prior():
166     return optimizer.apply_gradients(
167         grads_and_vars_except_prior,
168         global_step=global_step)
169
170 def train_op_all():
171     return optimizer.apply_gradients(
172         grads_and_vars,
173         global_step=global_step)
174
175 train_op = tf.cond(
176     global_step < params["prior_burn_in_steps"],
177     true_fn=train_op_except_prior,
178     false_fn=train_op_all)
179
180 # The perplexity is an exponent of the average negative ELBO per word.
181 words_per_document = tf.reduce_sum(features, axis=1)
182 log_perplexity = -elbo / words_per_document
183 tf.summary.scalar("perplexity", tf.exp(tf.reduce_mean(log_perplexity)))
184 (log_perplexity_tensor, log_perplexity_update) = tf.metrics.mean(
185     log_perplexity)
186 perplexity_tensor = tf.exp(log_perplexity_tensor)
187
188 # Obtain the topics summary. Implemented as a py_func for simplicity.
189 topics = tf.py_func(
190     functools.partial(get_topics_strings, vocabulary=params["vocabulary"]),
191     [topics_words, concentration], tf.string, stateful=False)
192 tf.summary.text("topics", topics)
193
194 return tf.estimator.EstimatorSpec(
195     mode=mode,
196     loss=loss,
197     train_op=train_op,
198     eval_metric_ops={
199         "elbo": tf.metrics.mean(elbo),
200         "log_likelihood": tf.metrics.mean(log_likelihood),
201         "kl": tf.metrics.mean(kl),
202         "perplexity": (perplexity_tensor, log_perplexity_update),
203         "topics": (topics, tf.no_op()),
204     },
205 )
206
207 def main(argv):
208     del argv # unused
209
210
```

```
214
215 if FLAGS.delete_existing and tf.gfile.Exists(FLAGS.model_dir):
216     tf.gfile.DeleteRecursively(FLAGS.model_dir)
217     tf.gfile.MakeDirs(FLAGS.model_dir)
218
219 if FLAGS.fake_data:
220     train_input_fn, eval_input_fn, vocabulary = build_fake_input_fns(
221         FLAGS.batch_size)
222 else:
223     train_input_fn, eval_input_fn, vocabulary = build_input_fns(
224         FLAGS.data_dir, FLAGS.batch_size)
225     params["vocabulary"] = vocabulary
226
227 estimator = tf.estimator.Estimator(
228     model_fn,
229     params=params,
230     config=tf.estimator.RunConfig(
231         model_dir=FLAGS.model_dir,
232         save_checkpoints_steps=FLAGS.viz_steps,
233     ),
234 )
235
236 for _ in range(FLAGS.max_steps // FLAGS.viz_steps):
237     estimator.train(train_input_fn, steps=FLAGS.viz_steps)
238     eval_results = estimator.evaluate(eval_input_fn)
239     # Print the evaluation results. The keys are strings specified in
240     # eval_metric_ops, and the values are Numpy scalars/arrays.
241     for key, value in eval_results.items():
242         print(key)
243         if key == "topics":
244             # Topics description is a np.array which prints better row-by-row.
245             for s in value:
246                 print(s)
247         else:
248             print(str(value))
249             print("")
250             print("")
251
252 if __name__ == "__main__":
253     tf.app.run()
254
```

Mesh TensorFlow

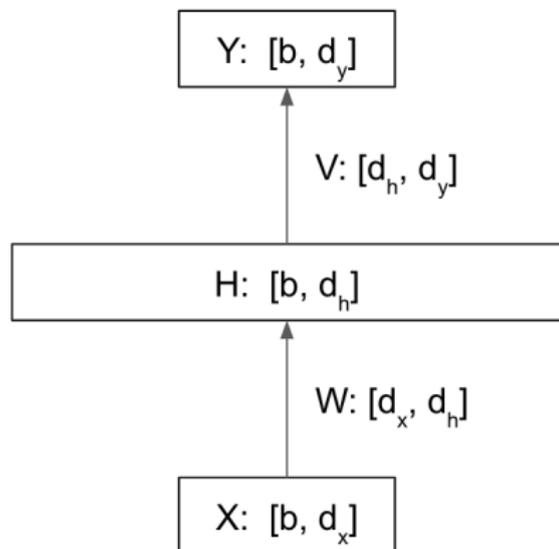
TPU Data Parallelism

- Parameters replicated on every core.
- Batch split between cores.
- Sum (allreduce) parameter gradients. (very efficient on locally-connected networks such as TPUs)

TPU Data Parallelism

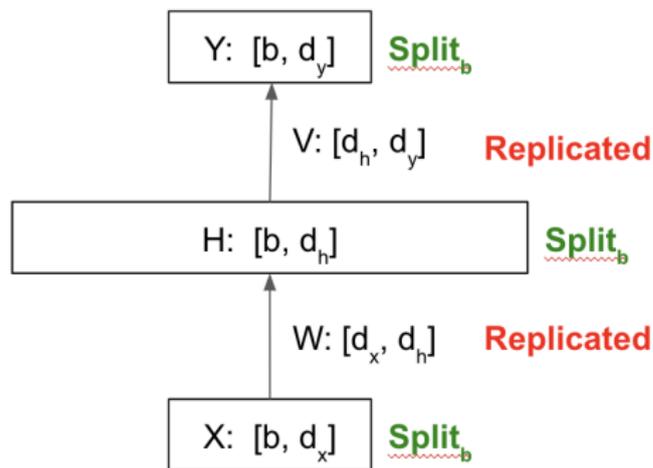
- Universal (any model/cluster)
- Fast to compile (SIMD)
- Full Utilization
- Allreduce is fast on any locally-connected network
- **All parameters must fit on one core.**

Example: Perceptron

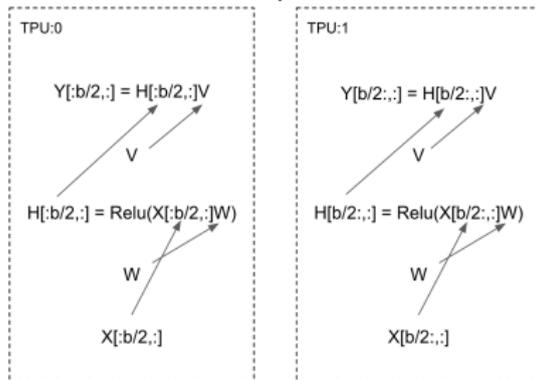


Example: Perceptron

$$Y = (H = \text{Relu}(XW_1))V$$

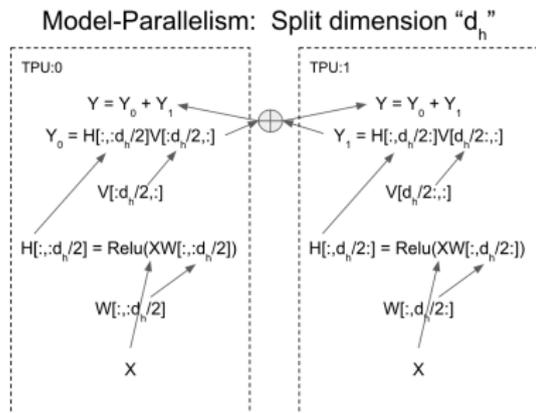
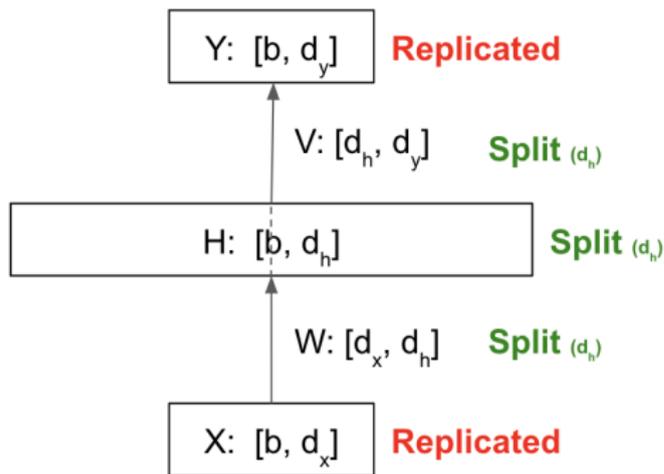


Data-Parallelism: Split dimension "b"



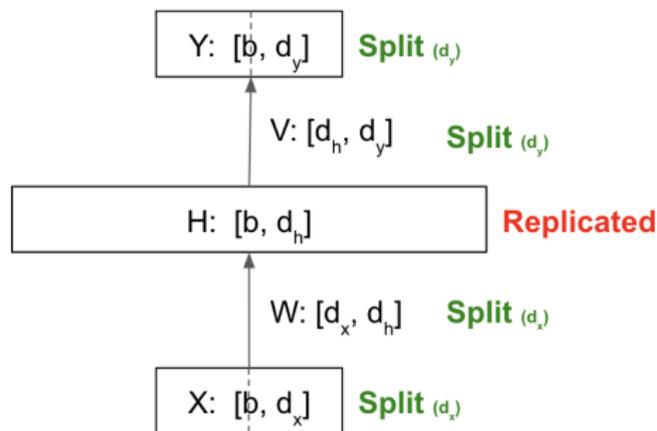
Example: Perceptron

$$Y = (H = \text{Relu}(XW_1))W_2$$

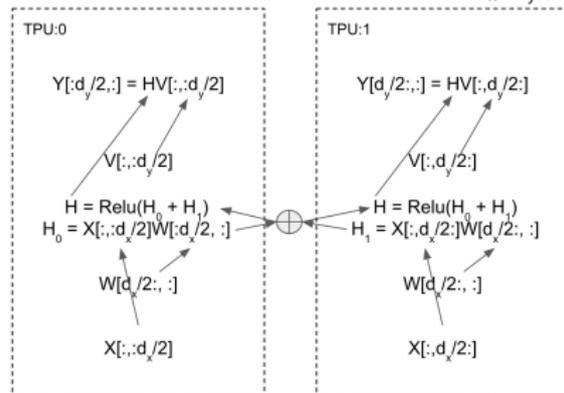


Example: Perceptron

$$Y = (H = \text{Relu}(XW_1))W_2$$



Model-Parallelism: Split dimensions d_x, d_y



Example: High-Quality Image Generation

50M+ parameter models (Image Transformer, VQVAE) on high-resolution images. Data parallelism.

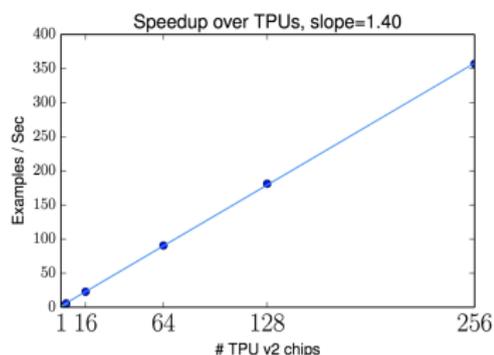


Figure 14: Vector-Quantized VAE on 64x64 ImageNet.

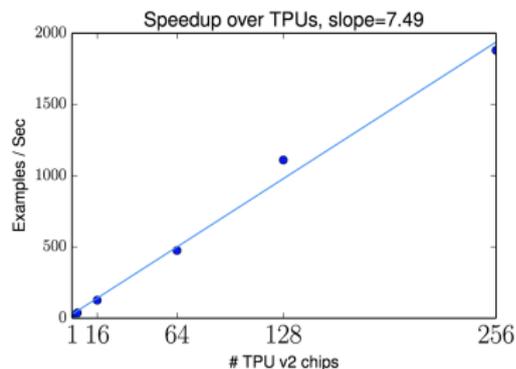


Figure 15: Image Transformer on 256x256 CelebA-HQ.

Edward2 achieves an optimal linear scaling from 1 to 256 TPUv2 chips.

Example: NUTS

Time per leapfrog step for No-U-Turn Sampler (NUTS) on Bayesian logistic regression. Covertypes, 500K data points, 54 features.

System	Runtime (ms)
Stan (CPU)	201.0
PyMC3 (CPU)	74.8
Handwritten TF (CPU)	66.2
Edward2 (CPU)	68.4
Handwritten TF (1 GPU)	9.5
Edward2 (1 GPU)	9.7
Edward2 (8 GPU)	2.3

Edward2 (GPU) achieves up to a 100x speedup over Stan and 7x over PyMC3. Dynamism is not possible in Edward 1.0.

Edward2 has negligible overhead over handwritten TF.

Example: Language Modeling

d_ff	heads	Parameters (in Billions)	Billion-Word Benchmark Word-Perplexity (bits/dim)	Wikipedia Subword-Perplexity (bits/dim)
8192	8	0.22	30.7	8.47
16384	16	0.37	28.0	7.92
32768	32	0.67	26.0	7.47
65516	64	1.28	24.6	6.97
131072	128	2.48	23.6	6.64
262144	256	4.90	23.1	6.41

Transformer from 20M to 3B parameter models. Model parallelism. Roughly 50% utilization.

Example: Machine Translation

d_ff	heads	Parameters (in Billions)	WMT'14 English-German BLEU	WMT'14 English-French BLEU
4096	8	0.24	28.6	41.7
8192	16	0.42	29.1	43.0
16384	32	0.77	28.9	43.2
32768	64	1.48	-	43.7
65516	128	2.89	-	43.7
4096	16	0.21	28.4	41.8 (Vaswani et. al)

Transformer from 20M to 3B parameter models. Model parallelism. Roughly 50% utilization.

Summary

1. Designing probabilistic systems for deep learning requires careful consideration about what's really brought to the table.
2. Our attempts pushed on what we think are the core elements.

Current directions.

1. We're advancing fundamental understandings of generative models and Bayesian neural networks.
2. We're pushing Mesh TensorFlow to trillion-parameter language models, new architectures, and model-parallel VAEs.

References

Systems

- Edward2: Simple, Distributed, Accelerated. NIPS 2018.
- Deep Learning for Supercomputers. NIPS 2018.
- Autoconj: Recognizing and Exploiting Conjugacy Without a Domain-Specific Language. NIPS 2018.

Methods

- Image Transformer. ICML 2018.
- Flipout: Efficient Pseudo-Independent Weight Perturbations on Mini-Batches. ICLR 2018.
- Reliable uncertainty estimates in deep neural networks using noise contrastive priors. arXiv:1807.09289 2018.



AutoConj: find and exploit exponential family structure without a DSL



Google AI

Matthew D. Hoffman*, Matthew J. Johnson*, Dustin V. Tran

TL;DR Write models in regular Python+Numpy with no mini-language, get exponential family structure-exploiting inference algorithms.

Why? Exploiting exponential family structure when it exists is labor-intensive, even for experts, which limits how we design new models and try new hybrid inference strategies (e.g. SVAEs). It's like neural nets before autodiff.

What is the autodiff for exponential family inference? **AutoConj!**

DSL? As with autodiff, don't want to be locked-in to a mini-language:

- New inference algorithms? Model classes?
- Optimization libraries? Automatic differentiation? Viz.?
- Compile to accelerators, distributed computing?

Need a system in native Python, and composable with others.

Background: exponential families

Define a probability model via a **statistic function** $t(x)$

$$p(x; \eta) = \exp\{(\eta, t(x)) - \mathcal{A}(\eta)\}, \quad \mathcal{A}(\eta) \triangleq \log \int \exp\{(\eta, t(x))\} \nu(dx),$$

Derivatives of the log partition function $\mathcal{A}(\eta)$ yield cumulants

$$\nabla \mathcal{A}(\eta) = \mathbb{E}[t(x)], \quad \nabla^2 \mathcal{A}(\eta) = \mathbb{E}[t(x)t(x)^T] - \mathbb{E}[t(x)] \mathbb{E}[t(x)]^T,$$

Compound models' statistics are **polynomials** in component statistics

$$\log p(z_1, z_2, \dots, z_M, x) = \sum_{i \in \mathcal{C}_A} (\eta_i(x), t_{i,1}(z_1)^{d_1} \otimes \dots \otimes t_{i,M}(z_M)^{d_M}) \\ \triangleq \mathcal{G}(\eta_{\mathcal{C}_A}(z_1, \dots, z_M), t_{\mathcal{C}_A}(x)),$$

Too much math for a poster

When g is multi-linear (has max-degree 1), then

Claim 2.1. Given an exponential family with density of the form (3), we have

$$p(z_{\mathcal{C}_A} | z_{-\mathcal{C}_A}) = \exp\{(\eta_{\mathcal{C}_A}^*(z_{-\mathcal{C}_A}), t_{\mathcal{C}_A}(z_{\mathcal{C}_A})) - \mathcal{A}_{\mathcal{C}_A}(\eta_{\mathcal{C}_A}^*)\} \text{ where } \eta_{\mathcal{C}_A}^* \triangleq \nabla_{\eta_{\mathcal{C}_A}} \mathcal{G}(\eta_{\mathcal{C}_A}(z_1, \dots, z_M), t_{\mathcal{C}_A}(x)).$$

Define a variational family using the same component statistics

$$q(x) = \prod_{i \in \mathcal{C}_A} q_i(x_i | \eta_{i, \mathcal{C}_A}), \quad q_i(x_i | \eta_{i, \mathcal{C}_A}) = \exp\{(\eta_{i, \mathcal{C}_A}(z_{-\mathcal{C}_A}), t_{i, \mathcal{C}_A}(x_i)) - \mathcal{A}_{i, \mathcal{C}_A}(\eta_{i, \mathcal{C}_A})\},$$

$$\log p(x) = \log \int p(x, z_{\mathcal{C}_A}) q(z_{\mathcal{C}_A}) dz_{\mathcal{C}_A} = \log \mathbb{E}_{q(z_{\mathcal{C}_A})} \left[\frac{p(x, z_{\mathcal{C}_A})}{q(z_{\mathcal{C}_A})} \right] \geq \mathbb{E}_{q(z_{\mathcal{C}_A})} \left[\log \frac{p(x, z_{\mathcal{C}_A})}{q(z_{\mathcal{C}_A})} \right] \triangleq \mathcal{L}.$$

Claim 2.2. Given a model with density of the form (3) and variational posterior (4), (5), we have $\arg \max_{\eta_{\mathcal{C}_A}} \mathcal{L}(\eta_{\mathcal{C}_A}) = \nabla_{\eta_{\mathcal{C}_A}} \mathcal{G}(\eta_{\mathcal{C}_A}, \dots, \eta_{\mathcal{C}_A})$ where $\eta_{i, \mathcal{C}_A} \triangleq \nabla_{\eta_{i, \mathcal{C}_A}} \mathcal{A}_{i, \mathcal{C}_A}(\eta_{i, \mathcal{C}_A})$, $\eta^i = 1, \dots, M$.

A general view on conjugacy: punchlines

- When energy is a **multi-linear polynomial in tractable statistic functions...**
 - Generic Gibbs via autodiff and a sampler for each statistic
 - Generic structured mean field and SVI via autodiff and a log normalizer for each statistic
 - Generic marginalization via autodiff and a log normalizer for each statistic
- Can write **generic implementations of structure-exploiting algorithms...**
 - but only once we've given the **polynomial representation**
 - ... and those are hard to write directly!
- Find polynomial representations automatically?

Term rewriting problem statement

Given a Python function denoting $f: \mathbb{R}^n \rightarrow \mathbb{R}$ that has a representation

$$f = g \circ h \quad \text{for a multi-linear polynomial } g: \mathbb{R}^{n_1} \times \dots \times \mathbb{R}^{n_M} \rightarrow \mathbb{R},$$

where the coordinate functions $h = (h_1, \dots, h_M)$ come from a known set,

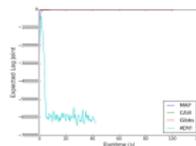
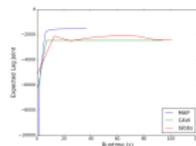
1. identify each h_i , and
2. produce a Python function to evaluate g .

Domain-specific term graph rewriting implementation

- **Tracer** using Autograd's API to map Python to term graphs
- **Pattern matcher** to do pattern-directed invocation
 - Python-embedded pattern language
 - Compiled into continuation-passing **matcher combinators** (~300 loc)
- **Rewriters** are syntactic graph macros using tracing to get **quasi-quasiquotes**

```

def rewrite(formula, op, x, y, arg1, arg2):
    return op(lap.elt(formula, <target> (x), <arg2>),
              op.elim(formula, <target> (y), <arg2>))
define_rewrite_rule('op', Add('op'), Val('x'), Segment('<arg2>'))
      
```



```

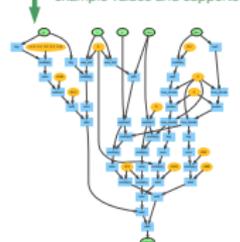
def normal_logpdf(x, loc, scale):
    prec = 1. / scale**2
    return -0.5 * prec * (x - loc)**2 - 0.5 * prec * log(prec)
           + 0.5 * prec * log(2 * pi)

def log_joint(pz, z, mu, tau, K):
    logp = 0.5 * sum((alpha - mu)**2) * log(prec)
           + 0.5 * sum(gamma * (alpha - mu))
    logp = normal_logpdf(mu, mu, tau) + log.agg(terms = tau)
    logp = op.sum(elt(formula, K) * op.log(pz))
    logp = ((-1)**op.log(tau)) * -0.5 * sum(log(pz))
           + gamma * (alpha)

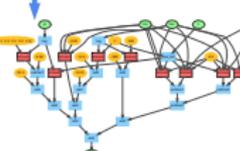
mu, s = op.diff(formula, K), mu
loglike = normal_logpdf(x, mu, s) * op.agg(terms)

return logp + loglike
      
```

1 Trace log joint density given example values and supports



2 Rewrite term graph to expose exponential family structure



3 Generic implementations of mean field, marginalization, Gibbs, etc. (in plain Python)

Model evaluation should be a first-class citizen in probabilistic programming

Alp Kucukelbir, Yixin Wang, Dustin Tran, David M. Blei
Columbia CS Columbia Stats Columbia CS Columbia CS + Stats
Fero Labs Google



1 | Introduction

Probabilistic programming research has been tightly focused on two things:

modeling and inference.

We argue that model evaluation deserves a similar level of attention.

Probabilistic programming enables the modern applied probabilist to craft bespoke probability models and perform inference with them. She can encode domain specific knowledge into her models with ease and express rich assumptions about the data she seeks to analyze. **With this freedom comes a pronounced need to evaluate such models.** Is there evidence for these assumptions? How well do these models work? We show how probabilistic programming languages offer practical solutions to some of these problems, but argue **that model evaluation deserves more interest from the community at large.**

2 | Methods for Model Evaluation

Focus | probability models with well-defined, evaluable joint distributions.

Scoring rules and point-wise evaluations

evaluating likelihood, computing losses, ideas around cross validation, posterior dispersion indices (Kucukelbir et al.).

Posterior predictive checks (PPCs)

1. Choose a statistic (e.g. min, max)
2. Simulate datasets from posterior predictive
3. Calculate statistics on simulated data
4. Compare to statistic evaluated on original data

Kernel-based methods

visualize smooth regions of data that is poorly explained by model (Lloyd and Ghahramani), kernel goodness-of-fit tests (Chwialkowski et al. | IJEU et al.)

3 | Status Quo and Future of Model Evaluation in Probabilistic Programming

Status Quo

Most popular probabilistic programming frameworks offer none or limited high-level constructs to implement model evaluation. Performing model evaluation in these cases requires manual implementation of the methods in Section 2.

Stan offers a helpful structure that aids in implementing model evaluation. For example, the generated quantities section can be used to compute PPCs and evaluate losses.

PyMC3 and Edward offer a productive out-of-the-box experience for model evaluation. Both have built-in implementations of PPCs and explicit documentation to do model evaluation and comparison. PyMC3 implements information criteria and Edward offers a suite of default scoring rules.

Future

The languages that facilitate model evaluation empower its users to build accurate and powerful probability models; this is a key goal for all probabilistic programming languages.

However, model evaluation faces its own set of challenges, unique to its application within probabilistic programming. Almost all automated inference algorithms are approximate. What happens to our evaluation metrics when the posterior is poor? Samples from MCMC algorithms may not have converged. Using a variational lower bound to the evidence can be dangerous for model comparison. PPCs may be incorrect due to approximation errors in the posterior distribution. Figure 1 shows an example of how this might occur.

Another open question is how to best integrate model evaluation into language semantics. Given the approximate nature of probabilistic programming inference algorithms, there are no accuracy guarantees for posterior computation under bounded time. How can language designers improve the language itself to expose the approximate nature of posterior computations and aid model evaluation?

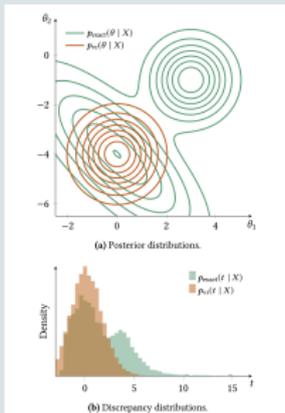


Figure 1. Assume a model with two latent variables $\theta = (\theta_1, \theta_2)$ and likelihood $\prod_{i=1}^N X_i(\theta_1, \theta_2, y_i, \sigma^2 = \exp(\theta_2))$. Sub-panel (a) depicts a bimodal distribution. A variational approximation to this posterior may capture only one of the modes. Sub-panel (b) shows the inaccuracy of the discrepancy distribution, as computed by a posterior predictive check, of $T = \max(X)$ using the variational posterior. Evaluating this model based on $p_{\text{var}}(T | X)$ could lead to incorrect conclusions.

Figure 1 shows an example of how this might occur.