# TensorFlow Distributions

Joshua V. Dillon[*], Ian Langmore[*], Dustin Tran[*†], Eugene Brevdo[*], Srinivas Vasudevan[*],
Dave Moore[*], Brian Patton[*], Alex Alemi[*], Matt Hoffman[*], Rif A. Saurous[*]

[*]Google, [†]Columbia University

## Abstract

The TensorFlow Distributions library implements a vision of probability theory adapted to the modern deep-learning paradigm of end-to-end differentiable computation. Building on two basic abstractions, it offers flexible building blocks for probabilistic computation. `Distributions` provide fast, numerically stable methods for generating samples and computing statistics, e.g., log density. `Bijectors` provide composable volume-tracking transformations with automatic caching. Together these enable modular construction of high dimensional distributions and transformations not possible with previous libraries (e.g., pixelCNNs, autoregressive flows, and reversible residual networks). They are the workhorse behind deep probabilistic programming systems like Edward and empower fast black-box inference in probabilistic models built on deep-network components. TensorFlow Distributions has proven an important part of the TensorFlow toolkit within Google and in the broader deep learning community.

***Keywords***   probabilistic programming, deep learning, probability distributions, transformations

## 1   Introduction

The success of deep learning—and in particular, deep generative models—presents exciting opportunities for probabilistic programming. Innovations with deep probabilistic models and inference algorithms have enabled new successes in perceptual domains such as images [19], text [4], and audio [48]; and they have advanced scientific applications such as understanding mouse behavior [23], learning causal factors in genomics [45], and synthesizing new drugs and materials [18].

Reifying these applications in code falls naturally under the scope of probabilistic programming systems, systems which build and manipulate computable probability distributions. Within the past year, languages for deep probabilistic programming such as Edward [46] have expanded deep-learning research by enabling new forms of experimentation, faster iteration cycles, and improved reproducibility.

```
e = make_encoder(x)
z = e.sample(n)
d = make_decoder(z)
r = make_prior()
avg_elbo_loss = tf.reduce_mean(
  e.log_prob(z) - d.log_prob(x) - r.log_prob(z))
train = tf.train.AdamOptimizer().minimize(
  avg_elbo_loss)
```

**Figure 1.** General pattern for training a variational auto-encoder (VAE) [27].

While there have been many developments in probabilistic programming languages, there has been limited progress in backend systems for probability distributions. This is despite their fundamental necessity for computing log-densities, sampling, and statistics, as well as for manipulations when composed as part of probabilistic programs. Existing distributions libraries lack modern tools necessary for deep probabilistic programming. Absent are: batching, automatic differentiation, GPU support, compiler optimization, composability with numerical operations and higher-level modules such as neural network layers, and efficient memory management.

To this end, we describe TensorFlow Distributions (r1.4), a TensorFlow (TF) library offering efficient, composable manipulation of probability distributions.[1]

**Illustration.** Figure 1 presents a template for a variational autoencoder under the TensorFlow Python API;[2] this is a generative model of binarized MNIST digits trained using amortized variational inference [27]. Figure 2 implements a standard version with a Bernoulli decoder, fully factorized Gaussian encoder, and Gaussian prior. By changing a few lines, Figure 3 implements a state-of-the-art architecture: a PixelCNN++ [41] decoder and a convolutional encoder and prior pushed through autoregressive flows [26, 36]. (`convnet`,`pixelcnnpp` are omitted for space.)

Figures 1 to 3 demonstrate the power of TensorFlow Distributions: fast, idiomatic modules are composed to express rich, deep structure. Section 5 demonstrates more applications: kernel density estimators, pixelCNN as a

---

[1]Home: tensorflow.org; Source: github.com/tensorflow/tensorflow.
[2]Namespaces: `tf=tensorflow`; `tfd=tf.contrib.distributions`; `tfb=tf.contrib.distributions.bijectors`.

1

```python
def make_encoder(x, z_size=8):
  net = make_nn(x, z_size*2)
  return tfd.MultivariateNormalDiag(
    loc=net[..., :z_size],
    scale=tf.nn.softplus(net[..., z_size:])))

def make_decoder(z, x_shape=(28, 28, 1)):
  net = make_nn(z, tf.reduce_prod(x_shape))
  logits = tf.reshape(
    net, tf.concat([[-1], x_shape], axis=0))
  return tfd.Independent(tfd.Bernoulli(logits))

def make_prior(z_size=8, dtype=tf.float32):
  return tfd.MultivariateNormalDiag(
    loc=tf.zeros(z_size, dtype)))

def make_nn(x, out_size, hidden_size=(128,64)):
  net = tf.flatten(x)
  for h in hidden_size:
    net = tf.layers.dense(
      net, h, activation=tf.nn.relu)
  return tf.layers.dense(net, out_size)
```

**Figure 2.** Standard VAE on MNIST with mean-field Gaussian encoder, Gaussian prior, Bernoulli decoder.

fully-visible likelihood, and how TensorFlow Distributions is used within higher-level abstractions (Edward and TF Estimator).

**Contributions.** TensorFlow Distributions (r1.4) defines two abstractions: `Distributions` and `Bijectors`. `Distributions` provides a collection of 56 distributions with fast, numerically stable methods for sampling, computing log densities, and many statistics. `Bijectors` provides a collection of 22 composable transformations with efficient volume-tracking and caching.

TensorFlow Distributions is integrated with the TensorFlow ecosystem [1]: for example, it is compatible with `tf.layers` for neural net architectures, `tf.data` for data pipelines, TF Serving for distributed computing, and TensorBoard for visualizations. As part of the ecosystem, TensorFlow Distributions inherits and maintains integration with TensorFlow graph operations, automatic differentiation, idiomatic batching and vectorization, device-specific kernel optimizations and XLA, and accelerator support for CPUs, GPUs, and tensor processing units (TPUs) [25].

TensorFlow Distributions is widely used in diverse applications. It is used by production systems within Google and by Google Brain and DeepMind for research prototypes. It is the backend for Edward [47].

## 1.1 Goals

TensorFlow Distributions is designed with three goals:

```python
import convnet, pixelcnnpp

def make_encoder(x, z_size=8):
  net = convnet(x, z_size*2)
  return make_arflow(
    tfd.MultivariateNormalDiag(
      loc=net[..., :z_size],
      scale_diag=net[..., z_size:])),
    invert=True)

def make_decoder(z, x_shape=(28, 28, 1)):
  def _logit_func(features):
    # implement single autoregressive step,
    # combining observed features with
    # conditioning information in z.
    cond = tf.layers.dense(z,
      tf.reduce_prod(x_shape))
    cond = tf.reshape(cond, features.shape)
    logits = pixelcnnpp(
      tf.concat((features, cond), -1))
    return logits
  logit_template = tf.make_template(
    "pixelcnn++", _logit_func)
  make_dist = lambda x: tfd.Independent(
    tfd.Bernoulli(logit_template(x)))
  return tfd.Autoregressive(
    make_dist, tf.reduce_prod(x_shape))

def make_prior(z_size=8, dtype=tf.float32):
  return make_arflow(
    tfd.MultivariateNormalDiag(
      loc=tf.zeros([z_size], dtype)))

def make_arflow(z_dist, n_flows=4,
  hidden_size=(640,)*3, invert=False):
  maybe_invert = tfb.Invert if invert else tfb.\
    Identity
  chain = list(itertools.chain.from_iterable([
    maybe_invert(tfb.MaskedAutoregressiveFlow(
      shift_and_log_scale_fn=tfb.\
      masked_autoregressive_default_template(
        hidden_size))),
    tfb.Permute(np.random.permutation(n_z)),
  ] for _ in range(n_flows)))
  return tfd.TransformedDistribution(
    distribution=z_dist,
    bijector=tfb.Chain(chain[:-1]))
```

**Figure 3.** State-of-the-art architecture. It uses a Pixel-CNN++ decoder [41] and autoregressive flows [26, 36] for encoder and prior.

**Fast.** TensorFlow Distributions is computationally and memory efficient. For example, it strives to use only

XLA-compatible ops (which enable compiler optimizations and portability to mobile devices), and whenever possible it uses differentiable ops (to enable end-to-end automatic differentiation). Random number generators for sampling call device-specific kernels implemented in C++. Functions with `Tensor` inputs also exploit vectorization through batches (Section 3.3). Multivariate distributions may be able to exploit additional vectorization structure.

**Numerically Stable.** All operations in TensorFlow Distributions are numerically stable across half, single, and double floating-point precisions (as TensorFlow `dtypes`: `tf.bfloat16` (truncated floating point), `tf.float16`, `tf.float32`, `tf.float64`). Class constructors have a `validate_args` flag for numerical asserts.

**Idiomatic.** As part of the TensorFlow ecosystem, TensorFlow Distributions maintains idioms such as inputs and outputs following a "`Tensor`-in, `Tensor`-out" pattern (though deviations exist; see Section 3.5), outputs preserving the inputs' `dtype`, and preferring statically determined shapes. Similarly, TensorFlow Distributions has no library dependencies besides NumPy [50] and six [37], further manages `dtypes`, supports TF-style broadcasting, and simplifies shape manipulation.

### 1.2 Non-Goals

TensorFlow Distributions does not cover all use-cases. Here we highlight goals common to probabilistic programming languages which are specifically not goals of this library.[3]

**Universality.** In order to be fast, the `Distribution` abstraction makes an explicit restriction on the class of computable distributions. Namely, any `Distribution` should offer `sample` and `log_prob` implementations that are computable in expected polynomial time. For example, the Multinomial-LogisticNormal distribution [7] fails to meet this contract.

This also precludes supporting a distributional calculus. For example, convolution is generally not analytic, so `Distributions` do not support the `__add__` operator: if $X \sim f_X$, $Y \sim f_Y$, and share domain $D$, then $Z = X + Y$ implies $f_Z(z) = \int_D f_X(z - y) f_Y(y) dy = (f_X * f_Y)(z)$.[4]

**Approximate Inference.** `Distributions` do not implement approximate inference or approximations of properties and statistics. For example, a Monte Carlo approximation of entropy is disallowed yet a function which computes an analytical, deterministic bound on

entropy is allowed. Compound distributions with conjugate priors such as Dirichlet-Multinomial are allowed. The marginal distribution of a hidden Markov model is also allowed since hidden states can be efficiently collapsed with the forward-backward algorithm [33].

## 2 Related Work

The R statistical computing language [21] provides a comprehensive collection of probability distributions. It inherits from the classic S language [6] and has accumulated user contributions over decades. We use R's collection as a goal for comprehensiveness and ease of user contribution. TensorFlow Distributions differs in being object-oriented instead of functional, enabling manipulation of `Distribution` objects; operations are also designed to be fast and differentiable. Most developers of the TensorFlow ecosystem are also Google-employed, meaning we benefit from more unification than R's ecosystem. For example, the popular `glmnet` and `lme4` R packages support only specific distributions for model-specific algorithms; all `Distributions` support generic TensorFlow optimizers for training/testing.

The SciPy `stats` module in Python collects probability distributions and statistical functions [24]. TensorFlow's primary demographic is machine learning users and researchers; they typically use Python. Subsequently, we modelled our API after SciPy; this mimics TensorFlow's API modelled after NumPy. Beyond API, the design details and implementations drastically differ. For example, TensorFlow Distributions enables arbitrary tensor-dimensional vectorization, builds operations in the TensorFlow computational graph, supports automatic differentiation, and can run on accelerators. The TensorFlow Distributions API also introduces innovations such as higher-order distributions (Section 3.5), distribution functionals (Section 3.6), and `Bijectors` (Section 4).

Stan Math [10] is a C++ templated library for numerical and statistical functions, and with automatic differentiation as the backend for the Stan probabilistic programming language [9]. Different from Stan, we focus on enabling deep probabilistic programming. This lead to new innovations with bijectors, shape semantics, higher-order distributions, and distribution functionals. Computationally, TensorFlow Distributions also enables a variety of non-CPU accelerators, and compiler optimizations in static over dynamically executed graphs.

## 3  Distributions

TensorFlow Distributions provides a collection of approximately 60 distributions with fast, numerically stable methods for sampling, log density, and many statistics. We describe key properties and innovations below.

---

[3]Users can subclass `Distribution` relaxing these properties.

[4]In future work, we may support this operation in cases when it satisfies our goals, e.g., for the analytic subset of stable distributions such as `Normal`, `Levy`.

## 3.1 Constructor

TensorFlow `Distributions` are object-oriented. A distribution implementation subclasses the `Distribution` base class. The base class follows a "public calls private" design pattern where, e.g., the public `sample` method calls a private `_sample` implemented by each subclass. This handles basic argument validation (e.g., type, shape) and simplifies sharing function documentation.

`Distributions` take the following arguments:

| | |
|---|---|
| **parameters** | indexes to family |
| **dtype** | dtype of samples |
| **reparameterization_type** | sampling (Section 3.4) |
| **validate_args** | bool permitting numerical checks |
| **allow_nan_stats** | bool permitting NaN outputs |
| **name** | str prepended to TF ops |

Parameter arguments support TF-style broadcasting. For example, `Normal(loc=1., scale=[0.5, 1., 1.5])` is effectively equivalent to `Normal(loc=[1., 1., 1.], scale=[0.5, 1., 1.5])`. Distributions use self-documenting argument names from a concise lexicon. We never use Greek and prefer, for example, `loc`, `scale`, `rate`, `concentration`, rather than $\mu$, $\sigma$, $\lambda$, $\alpha$.

Alternative parameterizations can sometimes lead to an "argument zoo." To migitate this, we distinguish between two cases. When numerical stability necessitates them, distributions permit mutually exclusive parameters (this produces only one extra argument). For example, `Bernoulli` accepts `logits` or `probs`, `Poisson` accepts `rate` or `log_rate`; neither permits specifying both. When alternative parameterizations are structural, we specify different classes: `MultivariateNormalTriL`, `MultivariateNormalDiag`, `MultivariateNormalDiagPlusLowRank` implement multivariate normal distributions with different covariance structures.

The `dtype` defaults to floats or ints, depending on the distribution's support, and with precision given by its parameters'. Distributions over integer-valued support (e.g., `Poisson`) use `tf.int*`. Distributions over real-valued support (e.g., `Dirichlet`) use `tf.float*`. This distinction exists because of mathematical consistency; and in practice, integer-valued distributions are often used as indexes into `Tensors`.[5]

If `validate_args=True`, argument validation happens during graph construction when possible; any validation at graph execution (runtime) is gated by a Boolean.

---

[5]Currently, TensorFlow Distributions' dtype does not follow this standard. For backwards compatibility, we are in the progress of implementing it by adding a new `sample_dtype` kwarg.

Among other checks, `validate_args=True` limits integer distributions' support to integers exactly representable by same-size IEEE754 floats, i.e., integers cannot exceed $2^{\text{significand\_bits}}$. If `allow_nan_stats=True`, operations can return NaNs; otherwise an error is raised.

## 3.2 Methods

At minimum, supported `Distributions` implement the following methods: `sample` to generate random outcome `Tensors`, `log_prob` to compute the natural logarithm of the probability density (or mass) function of random outcome `Tensors`, and `batch_shape_tensor`, `event_shape_tensor` to describe the dimensions of random outcome `Tensors` (Section 3.3), returned itself as `Tensors`.

Supported `Distributions` implement many additional methods, including `cdf`, `survival_function`, `quantile`, `mean`, `variance`, and `entropy`. The `Distribution` base class automates implementation of related functions such as `prob` given `log_prob` and `log_survival_fn` given `log_cdf` (unless a more efficient or numerically stable implementation is available). Distribution-specific statistics are permitted; for example, `Wishart` implements the expected log determinant (`mean_log_det`) of matrix variates, which would not be meaningful for univariate distributions.

All methods of supported distributions satisfy the following contract:

**Efficiently Computable.** All member functions have expected polynomial-time complexity. Further, they are vectorized (Section 3.3) and have optimized sampling routines (Section 3.4). TensorFlow Distributions also favors efficient parameterizations: for example, we favor `MultivariateNormalTriL`, whose covariance is parameterized by the outer product of a lower triangular matrix, over `MultivariateNormalFullCovariance` which requires a Cholesky factorization.

**Statistically Consistent.** Under `sample`, the Monte Carlo approximation of any statistic converges to the statistic's implementation as the number of samples approaches $\infty$. Similarly, `pdf` is equal to the derivative of `cdf` with respect to its input; and `sample` is equal in distribution to uniform sampling followed by the inverse of `cdf`.

**Analytical.** All member functions are analytical excluding `sample`, which is non-deterministic. For example, `Mixture` implements an analytic expression for an entropy lower bound method, `entropy_lower_bound`; its exact entropy is intractable. However, no method function's implementation uses a Monte Carlo estimate (even with a fixed seed, or low-discrepancy sequence [35]) which we qualify as non-analytical.

$$\begin{bmatrix} n \text{ Monte Carlo} \\ \text{draws} \end{bmatrix}, \begin{matrix} b \text{ examples per} \\ \text{batch} \end{matrix}, \begin{matrix} s \text{ latent} \\ \text{dimensions} \end{matrix}$$

<div style="text-align:center">

```
sample_shape        batch_shape         event_shape
(indep,             (indep, not         (can be
identically         identical)          dependent)
distributed)
```

</div>

**Figure 4.** Shape semantics. Refers to variational distribution in Figure 1.

**Fixed Properties.** In keeping with TensorFlow idioms, `Distribution` instances have fixed shape semantics (Section 3.3), `dtype`, class methods, and class properties throughout the instance's lifetime. Member functions have no side effects other than to add ops to the TensorFlow graph.

Note this is unlike the statefulness of exchangeable random primitives [2], where sampling can memoize over calls to lazily evaluate infinite-dimensional data structures. To handle such distributions, future work may involve a sampling method which returns another distribution storing those samples. This preserves immutability while enabling marginal representations of completely random measures such as a Chinese restaurant process, which is the compound distribution of a Dirichlet process and multinomial distribution [3]. Namely, its `sample` computes a Pólya urn-like scheme caching the number of customers at each table.[6]

### 3.3 Shape Semantics

To make `Distribution`s fast, a key challenge is to enable arbitrary tensor-dimensional vectorization. This allows users to properly utilize multi-threaded computation as well as array data-paths in modern accelerators. However, probability distributions involve a number of notions for different dimensions; they are often conflated and thus difficult to vectorize.

To solve this, we (conceptually) partition a `Tensor`'s shape into three groups:

1. *Sample shape* describes independent, identically distributed draws from the distribution.
2. *Batch shape* describes independent, *not* identically distributed draws. Namely, we may have a set of (different) parameterizations to the same distribution. This enables the common use case in machine learning of a "batch" of examples, each modelled by its own distribution.

---

[6]While the Chinese restaurant process is admittable as a (sequence of) `Distribution`, the Dirichlet process is not: its probability mass function involves a countable summation.

3. *Event shape* describes the shape of a single draw (event space) from the distribution; it may be dependent across dimensions.

Figure 4 illustrates this partition for the latent code in a variational autoencoder (Figure 1). Combining these three shapes in a single `Tensor` enables efficient, idiomatic vectorization and broadcasting.

Member functions all comply with the distribution's shape semantics and `dtype`. As another example, we initialize a batch of three multivariate normal distributions in $\mathbb{R}^2$. Each batch member has a different mean.

```
# Initialize 3-batch of 2-variate
# MultivariateNormals each with different mean.
mvn = tfd.MultivariateNormalDiag(
  loc=[[1., 1.], [2., 2.], [3., 3.]]))
x = mvn.sample(10)
# ==> x.shape=[10, 3, 2]. 10 samples across
#     3 batch members. Each sample in R^2.
pdf = mvn.prob(x)
# ==> pdf.shape=[10, 3]. One pdf calculation
#     for 10 samples across 3 batch members.
```

Partitioning `Tensor` dimensions by "sample", "batch", and "event" dramatically simplifies user code while naturally exploiting vectorization. For example, we describe a Monte Carlo approximation of a Normal-Laplace compound distribution,

$$p(x \mid \sigma, \mu_0, \sigma_0) = \int_{\mathbb{R}} \text{Normal}(x \mid \mu, \sigma) \, \text{Laplace}(\mu \mid \mu_0, \sigma_0) \, d\mu.$$

```
# Draw n iid samples from a Laplace.
mu = tfd.Laplace(
  loc=mu0, scale=sigma0).sample(n)
# ==> sample_shape = [n]
#     batch_shape = []
#     event_shape = []
# Compute n different Normal pdfs at
# scalar x, one for each Laplace draw.
pr_x_given_mu = tfd.Normal(
  loc=mu, scale=sigma).prob(x)
# ==> sample_shape = []
#     batch_shape = [n]
#     event_shape = []
# Average across each Normal's pdf.
pr_x = tf.reduce_mean(pr_x_given_mu, axis=0)
# ==> pr_estimate.shape=x.shape=[]
```

This procedure is automatically vectorized because the internal calculations are over tensors, where each represents the differently parameterized `Normal` distributions. `sigma` and `x` are automatically broadcast; their value is applied pointwise thus eliding `n` copies.

To determine batch and event shapes (sample shape is determined from the `sample` method), we perform

shape inference from parameters at construction time. Parameter dimensions beyond that necessary for a single distribution instance always determine batch shape. Inference of event shapes is typically not required as distributions often know it a priori; for example, `Normal` is univariate. On the other hand, `Multinomial` infers its event shape from the rightmost dimension of its `logits` argument. Dynamic sample and batch ranks are not allowed because they conflate the shape semantics (and thus efficient computation); dynamic event ranks are not allowed as a design choice for consistency.

Note that event shape (and shapes in general) reflects the numerical shape and not the mathematical definition of dimensionality. For example, `Categorical` has a scalar event shape over a finite set of integers; while a one-to-one mapping exists, `OneHotCategorical` has a vector event shape over one-hot vectors. Other distributions with non-scalar event shape include Dirichlet (simplexes) and Wishart (positive semi-definite matrices).

## 3.4 Sampling

Sampling is one of the most common applications of a `Distribution`. To optimize speed, we implement sampling by registering device-specific kernels in C++ to TensorFlow operations. We also use well-established algorithms for random number generation. For example, draws from `Normal` use the Box-Muller transform to return an independent pair of normal samples from an independent pair of uniform samples [8]; CPU, GPU, and TPU implementations exist. Draws from `Gamma` use the rejection sampling algorithm of Marsaglia and Tsang [30]; currently, only a CPU implementation exists.

**Reparameterization.** `Distributions` employ a `reparameterization_type` property (Section 3.1) which annotates the interaction between automatic differentiation and sampling. Currently, there are two such annotations: "fully reparameterized" and "not reparameterized". To illustrate "fully reparameterized", consider `dist = Normal(loc, scale)`. The sample `y = dist.sample()` is implemented internally via `x = tf.random_normal([]); y = scale * x + loc`. The sample y is "reparameterized" because it is a smooth function of the parameters `loc`, `scale`, and a parameter-free sample x. In contrast, the most common Gamma sampler is "not reparameterized": it uses an accept-reject scheme that makes the samples depend non-smoothly on the parameters [34].

When composed with other TensorFlow ops, a "fully reparameterized" `Distribution` enables end-to-end automatic differentiation on functions of its samples. A common use case is a loss depending on expectations of the form $E[\varphi(Y)]$ for some function $\varphi$. For example,

variational inference algorithms minimize the KL divergence between $p_Y$ and another distribution $h$, $E[\log[p_Y(Y)/h(Y)]]$ using gradient-based optimization. To do so, one can compute a Monte Carlo approximation

$$S_N := \frac{1}{N} \sum_{n=1}^{N} \varphi(Y_n), \text{ where } Y_n \sim p_Y. \qquad (1)$$

This lets us use $S_N$ not only as an estimate of our expected loss $E[\varphi(Y)]$, but also use $\nabla_\lambda S_N$ as an estimate of the gradient $\nabla_\lambda E[\varphi(Y)]$ with respect to parameters $\lambda$ of $p_Y$. If the samples $Y_n$ are reparameterized (in a smooth enough way), then both approximations are justified [14, 27, 42].

## 3.5 Higher-Order Distributions

Higher-order distributions are `Distributions` which are functions of other `Distributions`. This deviation from the `Tensor`-in, `Tensor`-out pattern enables modular, inherited construction of an enormous number of distributions. We outline three examples below and use a running illustration of composing distributions.

`TransformedDistribution` is a distribution $p(y)$ consisting of a base distribution $p(x)$ and invertible, differentiable transform $Y = g(X)$. The base distribution is an instance of the `Distribution` class and the transform is an instance of the `Bijector` class (Section 4).

For example, we can construct a (standard) Gumbel distribution from an Exponential distribution.

```
standard_gumbel = tfd.TransformedDistribution(
  distribution=tfd.Exponential(rate=1.),
  bijector=tfb.Chain([
    tfb.Affine(
      scale_identity_multiplier=-1.,
      event_ndims=0),
    tfb.Invert(tfb.Exp()),
  ]))
standard_gumbel.batch_shape # ==> []
standard_gumbel.event_shape # ==> []
```

The `Invert(Exp)` transforms the Exponential distribution by the natural-log, and the `Affine` negates. In general, algebraic relationships of random variables are powerful, enabling distributions to inherit method implementations from parents (e.g., internally, we implement multivariate normal distributions as `Affine` transforms of `Normal`).

Building on `standard_gumbel`, we can also construct $2*28*28$ independent relaxations of the Bernoulli distribution, known as Gumbel-Softmax or Concrete [22, 29].

```
alpha = tf.stack([
  tf.fill([28 * 28], 2.),
  tf.ones(28 * 28)])
```

```
concrete_pixel = tfd.TransformedDistribution(
  distribution=standard_gumbel,
  bijector=tfb.Chain([
    tfb.Sigmoid(),
    tfb.Affine(shift=tf.log(alpha)),
  ]),
  batch_shape=[2, 28 * 28])
concrete_pixel.batch_shape # ==> [2, 784]
concrete_pixel.event_shape # ==> []
```

The `Affine` shifts by `log(alpha)` for two batches. Applying `Sigmoid` renders a batch of $[2, 28*28]$ univariate Concrete distributions.

**Independent** enables idiomatic manipulations between batches and event dimensions. Given a `Distribution` instance `dist` with batch dimensions, Independent builds a vector (or matrix, or tensor) valued distribution whose event components default to the rightmost batch dimension of `dist`.

Building on `concrete_pixel`, we reinterpret the 784 batches as jointly characterizing a distribution.

```
image_dist = tfd.TransformedDistribution(
  distribution=tfd.Independent(concrete_pixel),
  bijector=tfb.Reshape(
    event_shape_out=[28, 28, 1],
    event_shape_in=[28 * 28]))
image_dist.batch_shape # ==> [2]
image_dist.event_shape # ==> [28, 28, 1]
```

The `image_dist` distribution is over $28 \times 28 \times 1$-dim. events (e.g., MNIST-resolution pixel images).

**Mixture** defines a probability mass function $p(x) = \sum_{k=1}^{K} \pi_k p(x \mid k)$, where the mixing weights $\pi_k$ are provided by a `Categorical` distribution as input, and the components $p(x \mid k)$ are arbitrary `Distributions` with same support. For components that are batches of the same family, `MixtureSameFamily` simplifies the construction where its components are from the rightmost batch dimension. Building on `image_dist`:

```
image_mixture = tfd.MixtureSameFamily(
  mixture_distribution=tfd.Categorical(
    probs=[0.2, 0.8]),
  components_distribution=image_dist)
image_mixture.batch_shape # ==> []
image_mixture.event_shape # ==> [28, 28, 1]
```

Here, `MixtureSameFamily` creates a mixture of two components with weights $[0.2, 0.8]$. The components are slices along the batch axis of `image_dist`.

### 3.6 Distribution Functionals

Functionals which take probability distribution(s) as input and return a scalar are ubiquitous. They include information measures such as entropy, cross entropy, and mutual information [11]; divergence measures such as Kullback-Leibler, Csiszár-Morimoto $f$-divergence [12, 31], and multi-distribution divergences [15]; and distance metrics such as integral probability metrics [32].

We implement all (analytic) distribution functionals as methods in `Distributions`. For example, below we write functionals of Normal distributions:

```
p = tfd.Normal(loc=0., scale=1.)
q = tfd.Normal(loc=-1., scale=2.)
xent = p.cross_entropy(q)
kl = p.kl_divergence(q)
# ==> xent - p.entropy()
```

## 4 Bijectors

We described `Distributions`, sources of stochasticity which collect properties of probability distributions. `Bijectors` are deterministic transformations of random outcomes and of equal importance in the library. It consists of 22 composable transformations for manipulating `Distributions`, with efficient volume-tracking and caching of pre-transformed samples. We describe key properties and innovations below.

### 4.1 Motivation

The `Bijector` abstraction is motivated by two challenges for enabling efficient, composable manipulations of probability distributions:

1. We seek a minimally invasive interface for manipulating distributions. Implementing every transformation of every distribution results in a combinatorial blowup and is not realistically maintainable. Such a policy is unlikely to keep pace with the pace of research. Lack of encapsulation exacerbates already complex ideas/code and discourages community contributions.
2. In deep learning, rich high-dimensional densities typically use invertible volume-preserving mappings or mappings with fast volume adjustments (namely, the logarithm of the Jacobian's determinant has linear complexity with respect to dimensionality) [36]. We'd like to efficiently and idiomatically support them.

By isolating stochasticity from determinism, `Distributions` are easy to design, implement, and validate. As we illustrate with the flexibility of `TransformedDistribution` in Section 3.5, the ability to simply swap out functions applied to the distribution is a surprisingly powerful asset. Programmatically, the `Bijector` distinction enables encapsulation and modular distribution constructions with inherited, fast method implementations. Statistically, `Bijectors` enable exploiting algebraic relationships among random variables.

## 4.2 Definition

To address Section 4.1, the `Bijector` API provides an interface for transformations of distributions suitable for any differentiable and bijective map (*diffeomorphism*) as well as certain non-injective maps (Section 4.5).

Formally, given a random variable $X$ and diffeomorphism $F$, we can define a new random variable $Y$ whose density can be written in terms of $X$'s,

$$p_Y(y) = p_X(F^{-1}(y)) \, |DF^{-1}(y)|, \qquad (2)$$

where $DF^{-1}$ is the inverse of the Jacobian of $F$. Each `Bijector` subclass corresponds to such a function $F$, and `TransformedDistribution` uses the bijector to automate the details of the transformation $Y = F(X)$'s density (Equation 2). This allows us to define *many* new distributions in terms of existing ones.

A `Bijector` implements how to transform a `Tensor` and how the input `Tensor`'s shape changes; this `Tensor` is presumed to be a random outcome possessing `Distribution` shape semantics. Three functions characterize how the `Tensor` is transformed:

1. `forward` implements $x \mapsto F(x)$, and is used by `TransformedDistribution.sample` to convert one random outcome into another. It also establishes the name of the bijector.
2. `inverse` undoes the transformation $y \mapsto F^{-1}(y)$ and is used by `TransformedDistribution.log_prob`.
3. `inverse_log_det_jacobian` computes $\log|DF^{-1}(y)|$ and is used by `TransformedDistribution.log_prob` to adjust for how the volume changes by the transformation. In certain settings, it is more numerically stable (or generally preferable) to implement the `forward_log_det_jacobian`. Because forward and reverse log ∘ det ∘ Jacobians differ only in sign, either (or both) may be implemented.

A `Bijector` also describes how it changes the `Tensor`'s shape so that `TransformedDistribution` can implement functions that compute event and batch shapes. Most bijectors do not change the `Tensor`'s shape. Those which do implement `forward_event_shape_tensor` and `inverse_event_shape_tensor`. Each takes an input shape (vector) and returns a new shape representing the `Tensor`'s event/batch shape after `forward` or `inverse` transformations. Excluding higher-order bijectors, currently only `SoftmaxCentered` changes the shape.[7]

Using a `Bijector`, `TransformedDistribution` automatically and efficiently implements `sample`, `log_prob`,

and `prob`. For bijectors with constant Jacobian such as `Affine`, `TransformedDistribution` automatically implements statistics such as `mean`, `variance`, and `entropy`. The following example implements an affine-transformed Laplace distribution.

```
vector_laplace = tfd.TransformedDistribution(
  distribution=tfd.Laplace(loc=0., scale=1.),
  bijector=tfb.Affine(
    shift=tf.Variable(tf.zeros(d)),
    scale_tril=tfd.fill_triangular(
      tf.Variable(tf.ones(d*(d+1)/2)))),
  event_shape=[d])
```

The distribution is learnable via `tf.Variable`s and that the `Affine` is parameterized by what is essentially the Cholesky of the covariance matrix. This makes the multivariate construction computationally efficient and more numerically stable; bijector caching (Section 4.4) may even eliminate back substitution.

## 4.3 Composability

`Bijector`s can compose using higher-order `Bijector`s such as `Chain` and `Invert`. Figure 3 illustrates a powerful example where the `arflow` method composes a sequence of autoregressive and permutation `Bijector`s to compactly describe an autoregressive flow [26, 36].

The `Chain` bijector enables simple construction of rich `Distributions`. Below we construct a multivariate logit-Normal with matrix-shaped events.

```
matrix_logit_mvn =
  tfd.TransformedDistribution(
    distribution=tfd.Normal(0., 1.),
    bijector=tfb.Chain([
      tfb.Reshape([d, d]),
      tfb.SoftmaxCentered(),
      tfb.Affine(scale_diag=diag),
    ]),
    event_shape=[d * d])
```

The `Invert` bijector effectively doubles the number of bijectors by swapping `forward` and `inverse`. It is useful in situations such as the Gumbel construction in Section 3.5. It is also useful for transforming constrained continuous distributions onto an unconstrained real-valued space. For example:

```
softminus_gamma = tfd.TransformedDistribution(
 distribution=tfd.Gamma(
   concentration=alpha,
   rate=beta),
 bijector=tfb.Invert(tfb.Softplus()))
```

---

[7]To implement $\text{softmax}(x) = \exp(x)/\sum_i \exp(x_i)$ as a diffeomorphism, its `forward` appends a zero to the event and its `reverse` strips this padding. The result is a bijective map which avoids the fact that $\text{softmax}(x) = \exp(x-c)/\sum_i \exp(x_i - c)$ for any $c$.

This performs a softplus-inversion to robustly transform Gamma to be unconstrained. This enables a key component of automated algorithms such as automatic differentiation variational inference [28] and the No U-Turn Sampler [20]. They only operate on real-valued spaces, so unconstraints expand their scope.

## 4.4 Caching

Bijectors automatically cache input/output pairs of operations, including the log ∘ det ∘ Jacobian. This is advantageous when the inverse calculation is slow, numerically unstable, or not easily implementable. A cache hit occurs when computing the probability of results of sample. That is, if $q(x)$ is the distribution associated with $x = f(\varepsilon)$ and $\varepsilon \sim r$, then caching lowers the cost of computing $q(x_i)$ since

$$q(x_i) = r((f^{-1} \circ f)(\varepsilon_i)) \left| \left( \frac{\partial}{\partial \varepsilon} \circ f^{-1} \circ f \right)(\varepsilon_i) \right|^{-1} = r(\varepsilon_i).$$

Because TensorFlow follows a deferred execution model, Bijector caching is nominal; it has zero memory or computational cost. The Bijector base class merely replaces one graph node for another already existing node. Since the existing node ("cache hit") is already an execution dependency, the only cost of "caching" is during graph construction.

Caching is computationally and numerically beneficial for importance sampling algorithms, which compute expectations. They weight by a drawn samples' reciprocal probability. Namely,

$$\begin{aligned} \mu &= \int f(x) p(x) dx \\ &= \int \frac{f(x) p(x)}{q(x)} q(x) dx \\ &= \lim_{n \to \infty} n^{-1} \sum_i^n \frac{f(x_i) p(x_i)}{q(x_i)}, \text{ where } x_i \overset{\text{iid}}{\sim} q. \end{aligned}$$

Caching also nicely complements black-box variational inference algorithms [28, 40]. In these procedures, the approximate posterior distribution only computes log_prob over its own samples. In this setting the sample's preimage ($\varepsilon_i$) is known without computing $f^{-1}(x_i)$.

MultivariateNormalTriL is implemented as a TransformedDistribution with the Affine bijector. Caching removes the need for quadratic complexity back substitution. For an InverseAutoregressiveFlow [26],

```
laplace_iaf = tfd.TransformedDistribution(
  distribution=tfd.Laplace(loc=0., scale=1.),
  bijector=tfb.Invert(
    tfb.MaskedAutoregressiveFlow(
      shift_and_log_scale_fn=tfb.\
      masked_autoregressive_default_template(
        hidden_layers))),
  event_shape=[d])
```

caching reduces the overall complexity from quadratic to linear (in event size).

## 4.5 Smooth Coverings

The Bijector framework extends to non-injective transformations, i.e., *smooth coverings* [44].[8] Formally, a smooth covering is a continuous function $F$ on the entire domain $D$ where, ignoring sets of measure zero, the domain can be partitioned as a finite union $D = D_1 \cup \cdots \cup D_K$ such that each restriction $F : D_k \to F(D)$ is a diffeomorphism. Examples include AbsValue and Square. We implement them by having the inverse method return the set inverse $F^{-1}(y) := \{x : F(x) = y\}$ as a tuple.

Smooth covering Bijectors let us easily build half-distributions, which allocate probability mass over only the positive half-plane of the original distribution. For example, we build a half-Cauchy distribution as

```
half_cauchy = tfd.TransformedDistribution(
  bijector=tfb.AbsValue(),
  distribution=tfd.Cauchy(loc=0., scale=1.))
```

The half-Cauchy and half-Student t distributions are often used as "weakly informative" priors, which exhibit heavy tails, for variance parameters in hierarchical models [16].

## 5 Applications

We described two abstractions: Distributions and Bijectors. Recall Figures 1 to 3, where we showed the power of combining these abstractions for changing from simple to state-of-the-art variational auto-encoders. Below we show additional applications of TensorFlow Distributions as part of the TensorFlow ecosystem.

### 5.1 Kernel Density Estimators

A kernel density estimator (KDE) is a nonparametric estimator of an unknown probability distribution [51]. Kernel density estimation provides a fundamental smoothing operation on finite samples that is useful across many applications. With TensorFlow Distributions, KDEs can be flexibly constructed as a MixtureSameFamily. Given a finite set of points x in $\mathbb{R}^D$, we write

```
f = lambda x: tfd.Independent(tfd.Normal(
  loc=x, scale=1.))
n = x.shape[0].value
kde = tfd.MixtureSameFamily(
  mixture_distribution=tfd.Categorical(
    probs=[1 / n] * n),
  components_distribution=f(x))
```

Here, f is a callable taking x as input and returns a distribution. Above, we use an independent $D$-dimensional

---

[8]Bijector caching is currently not supported for smooth coverings.

```
import pixelcnn

def pixelcnn_dist(params, x_shape=(32, 32, 3)):
  def _logit_func(features):
    # implement single autoregressive step
    # on observed features
    logits = pixelcnn(features)
    return logits
  logit_template = tf.make_template(
    "pixelcnn", _logit_func)
  make_dist = lambda x: tfd.Independent(
    tfd.Bernoulli(logit_template(x)))
  return tfd.Autoregressive(
    make_dist, tf.reduce_prod(x_shape)))

x = pixelcnn_dist()
loss = -tf.reduce_sum(x.log_prob(images))
train = tf.train.AdamOptimizer(
    ).minimize(loss)  # run for training
generate = x.sample()  # run for generation
```

**Figure 5.** PixelCNN distribution on images. It uses `Autoregressive`, which takes as input a callable returning a distribution per time step.

`Normal` distribution (equivalent to `MultivariateNormalDiag`), which induces a Gaussian kernel density estimator.

Changing the callable extends the KDE to alternative distribution-based kernels. For example, we can use `lambda x: MultivariateNormalTriL(loc=x)` for a multivariate kernel, and alternative distributions such as `lambda x: Independent(StudentT(loc=x, scale=0.5, df=3)`. The same concept also applies for bootstrap techniques [13]. We can employ parametric bootstrap or stratified sampling to replace the equal mixing weights.

## 5.2 PixelCNN Distribution

Building from the KDE example above, we now show a modern, high-dimensional density estimator. Figure 5 builds a PixelCNN [49] as a fully-visible autoregressive distribution on `images`, which is a batch of $32 \times 32 \times 3$ RGB pixel images from Small ImageNet.

The variable x is the pixelCNN distribution. It makes use of the higher-order `Autoregressive` distribution, which takes as input a Python callable and number of autoregressive steps. The Python callable takes in currently observed features and returns the per-time step distribution. `-tf.reduce_sum(x.log_prob(images))` is the loss function for maximum likelihood training; `x.sample` generates new images.

We also emphasize modularity. Note here, we used the pixelCNN as a fully-visible distribution. This differs from Figure 3 which employs pixelCNN as a decoder (conditional likelihood).

```
from edward.models import Normal

z = x = []
z[0] = Normal(loc=tf.zeros(K),scale=tf.ones(K))
h = tf.layers.dense(
  z[0], 512, activation=tf.nn.relu)
loc = tf.layers.dense(h, D, activation=None)
x[0] = Normal(loc=loc, scale=0.5)
for t in range(1, T):
  inputs = tf.concat([z[t - 1], x[t - 1]], 0)
  loc = tf.layers.dense(
    inputs, K, activation=tf.tanh)
  z[t] = Normal(loc=loc, scale=0.1)
  h = tf.layers.dense(
    z[t], 512, activation=tf.nn.relu)
  loc = tf.layers.dense(h, D, activation=None)
  x[t] = Normal(loc=loc, scale=0.5)
```

**Figure 6.** Stochastic recurrent neural network, which is a state space model with nonlinear dynamics. The transition mimics a recurrent tanh cell and the omission is multi-layered.

## 5.3 Edward Probabilistic Programs

We describe how TensorFlow Distributions enables Edward as a backend. In particular, note that non-goals in TensorFlow Distributions can be accomplished at higher-level abstractions. Here, Edward wraps TensorFlow Distributions as random variables, associating each `Distribution` to a random outcome `Tensor` (calling `sample`) in the TensorFlow graph. This enables a calculus where TensorFlow ops can be applied directly to Edward random variables; this is a non-goal for TensorFlow Distributions.

As an example, Figure 6 implements a stochastic recurrent neural network (RNN), which is an RNN whose hidden state is random [5]. Formally, for $T$ time steps, the model specifies the joint distribution

$$p(\mathbf{x}, \mathbf{z}) = \text{Normal}(\mathbf{z}_1 \mid \mathbf{0}, \mathbf{I}) \prod_{t=2}^{T} \text{Normal}(\mathbf{z}_t \mid f(\mathbf{z}_{t-1}), 0.1)$$

$$\prod_{t=1}^{T} \text{Normal}(\mathbf{x}_t \mid g(\mathbf{z}_t), 0.5),$$

where each time step in an observed real-valued sequence $\mathbf{x} = [\mathbf{x}_1, \dots, \mathbf{x}_T] \in \mathbb{R}^{T \times D}$ is associated with an unobserved state $\mathbf{z}_t \in \mathbb{R}^K$; the initial latent variable $\mathbf{z}_1$ is drawn randomly from a standard normal. The noise standard deviations are fixed and broadcasted over the batch. The latent variable and likelihood are parameterized by neural networks.

The program is generative: starting from a latent state, it unrolls state dynamics through time. Given this program and data, Edward's algorithms enable approximate

inference (a second non-goal of TensorFlow Distributions).

## 5.4 TensorFlow Estimator API

As part of the TensorFlow ecosystem, TensorFlow Distributions complements other TensorFlow libraries. We show how it complements TensorFlow Estimator.

Figure 7 demonstrates multivariate linear regression in the presence of heteroscedastic noise,

$$U \sim \text{MultivariateNormal}(0, I_d)$$

$$Y = \Sigma^{\frac{1}{2}}(X)U + \mu(X)$$

where $\Sigma : \mathbb{R}^d \to \{Z \in \mathbb{R}^{d \times d} : Z \succcurlyeq 0\}, \mu : \mathbb{R}^d \to \mathbb{R}^d$, and $\Sigma^{\frac{1}{2}}$ denotes the Cholesky decomposition. Adding more `tf.layers` to the parameterization of the `MultivariateNormalTriL` enables learning nonlinear transformations. ($\Sigma = I_d$ would be appropriate in homoscedastic regimes.)

Using `Distributions` to build `Estimators` is natural and ergonomic. We use `TPUEstimator` in particular, which extends `Estimator` with configurations for TPUs [25]. Together, `Distributions` and `Estimators` provide a simple, scalable platform for efficiently deploying training, evaluation, and prediction over diverse hardware and network topologies.

Figure 7 only writes the `Estimator` object. For training, call `estimator.train()`; for evaluation, call `estimator.evaluate()`; for prediction, call `estimator.predict()`. Each takes an input function to load in data.

## 6 Discussion

The TensorFlow Distributions library implements a vision of probability theory adapted to the modern deep-learning paradigm of end-to-end differentiable computation. `Distributions` provides a collection of 56 distributions with fast, numerically stable methods for sampling, computing log densities, and many statistics. `Bijectors` provides a collection of 22 composable transformations with efficient volume-tracking and caching.

Although Tensorflow Distributions is relatively new, it has already seen widespread adoption both inside and outside of Google. External developers have built on it to design probabilistic programming and statistical systems including Edward [47] and Greta [17]. Further, `Distribution` and `Bijector` is being used as the design basis for similar functionality in the PyTorch computational graph framework [39], as well as the Pyro and ZhuSuan probabilistic programming systems [38, 43].

Looking forward, we plan to continue expanding the set of supported `Distributions` and `Bijectors`. We

```python
def mvn_regression_fn(
    features, labels, mode, params=None):
  d = features.shape[-1].value
  mvn = tfd.MultivariateNormalTriL(
    loc=tf.layers.dense(features, d),
    scale_tril=tfd.fill_triangular(
      tf.layers.dense(features, d*(d+1)/2)))
  if mode == tf.estimator.ModeKeys.PREDICT:
    return mvn.mean()
  loss = -tf.reduce_sum(mvn.log_prob(labels))
  if mode == tf.estimator.ModeKeys.EVAL:
    metric_fn = lambda x,y:
      tf.metrics.mean_squared_error(x, y)
    return tpu_estimator.TPUEstimatorSpec(
      mode=mode,
      loss=loss,
      eval_metrics=(
        metric_fn, [labels, mvn.mean()]))
  optimizer = tf.train.AdamOptimizer()
  if FLAGS.use_tpu:
    optimizer = tpu_optimizer.
    CrossShardOptimizer(optimizer)
  train_op = optimizer.minimize(loss)
  return tpu_estimator.TPUEstimatorSpec(
    mode=mode, loss=loss, train_op=train_op)

# TPUEstimator Boilerplate.
run_config = tpu_config.RunConfig(
  master=FLAGS.master,
  model_dir=FLAGS.model_dir,
  session_config=tf.ConfigProto(
    allow_soft_placement=True,
    log_device_placement=True),
  tpu_config=tpu_config.TPUConfig(
    FLAGS.iterations,
    FLAGS.num_shards))
estimator = tpu_estimator.TPUEstimator(
  model_fn=mvn_regression_fn,
  config=run_config,
  use_tpu=FLAGS.use_tpu,
  train_batch_size=FLAGS.batch_size,
  eval_batch_size=FLAGS.batch_size)
```

**Figure 7.** Multivariate regression with TPUs.

intend to expand the `Distribution` interface to include supports, e.g., real-valued, positive, unit interval, etc., as a class property. We are also exploring the possibility of exposing exponential family structure, for example providing separate `unnormalized_log_prob` and `log_normalizer` methods where appropriate.

As part of the trend towards hardware-accelerated linear algebra, we are working to ensure that all distribution and bijector methods are compatible with TPUs [25], including special functions such as gamma, as well as rejection sampling-based (e.g., Gamma) and while-loop based sampling mechanisms (e.g., Poisson). We also aim to natively support Distributions over `SparseTensors`.

## Acknowledgments

## References

[1] Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow, Andrew Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Jozefowicz, Lukasz Kaiser, Manjunath Kudlur, Josh Levenberg, Dan Mané, Rajat Monga, Sherry Moore, Derek Murray, Chris Olah, Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda Viégas, Oriol Vinyals, Pete Warden, Martin Wattenberg, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. 2015. TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems. (2015). https://www.tensorflow.org/ Software available from tensorflow.org.

[2] Nathanael L Ackerman, Cameron E Freer, and Daniel M Roy. 2016. Exchangeable random primitives. In *Workshop on Probabilistic Programming Semantics*. 2016.

[3] David J Aldous. 1985. Exchangeability and related topics. In *École d'Été de Probabilités de Saint-Flour XIII—1983*. Springer, 1–198.

[4] Anonymous. 2017. Generative Models for Data Efficiency and Alignment in Language. *OpenReview* (2017).

[5] Justin Bayer and Christian Osendorfer. 2014. Learning Stochastic Recurrent Networks. *arXiv.org* (2014). arXiv:1411.7610v3

[6] Richard A Becker and John M Chambers. 1984. *S: an interactive environment for data analysis and graphics*. CRC Press.

[7] David M Blei and John Lafferty. 2006. Correlated topic models. In *Neural Information Processing Systems*.

[8] G. E. P. Box and Mervin E. Muller. 1958. A Note on the Generation of Random Normal Deviates. *The Annals of Mathematical Statistics* (1958), 610–611.

[9] Bob Carpenter, Andrew Gelman, Matthew D Hoffman, Daniel Lee, Ben Goodrich, Michael Betancourt, Marcus Brubaker, Jiqiang Guo, Peter Li, and Allen Riddell. 2016. Stan: a probabilistic programming language. *Journal of Statistical Software* (2016).

[10] Bob Carpenter, Matthew D Hoffman, Marcus Brubaker, Daniel Lee, Peter Li, and Michael Betancourt. 2015. The Stan math library: Reverse-mode automatic differentiation in C++. *arXiv preprint arXiv:1509.07164* (2015).

[11] Thomas M Cover and Joy A Thomas. 1991. *Elements of Information Theory*. Wiley Series in Telecommunications and Signal Processing.

[12] Imre Csiszár. 1963. Eine informationstheoretische Ungleichung und ihre Anwendung auf Beweis der Ergodizitaet von Markoffschen Ketten. *Magyer Tud. Akad. Mat. Kutato Int. Koezl.* 8 (1963), 85–108.

[13] Bradley Efron and Robert J Tibshirani. 1994. *An introduction to the bootstrap*. CRC press.

[14] M.C. Fu. 2006. *Simulation*. Handbook in Operations Research and Management Science, Vol. 13. North Holland.

[15] Dario Garcia-Garcia and Robert C Williamson. 2012. Divergences and risks for multiclass experiments. In *Conference on Learning Theory*.

[16] Andrew Gelman et al. 2006. Prior distributions for variance parameters in hierarchical models (comment on article by Browne and Draper). *Bayesian analysis* 1, 3 (2006), 515–534.

[17] Nick Golding. 2017. *greta: Simple and Scalable Statistical Modelling in R*. https://CRAN.R-project.org/package=greta R package version 0.2.0.

[18] Rafael Gómez-Bombarelli, David Duvenaud, José Miguel Hernández-Lobato, Jorge Aguilera-Iparraguirre, Timothy D Hirzel, Ryan P Adams, and Alán Aspuru-Guzik. 2016. Automatic chemical design using a data-driven continuous representation of molecules. *arXiv preprint arXiv:1610.02415* (2016).

[19] Ian Goodfellow, Jean Pouget-Abadie, M Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2014. Generative Adversarial Nets. In *Neural Information Processing Systems*.

[20] Matthew D Hoffman and Andrew Gelman. 2014. The no-U-turn sampler: adaptively setting path lengths in Hamiltonian Monte Carlo. *Journal of Machine Learning Research* 15 (2014), 1593–1623.

[21] Ross Ihaka and Robert Gentleman. 1996. R: A Language for Data Analysis and Graphics. *Journal of Computational and Graphical Statistics* 5, 3 (1996), 299–314. https://doi.org/10.1080/10618600.1996.10474713

[22] Eric Jang, Shixiang Gu, and Ben Poole. 2017. Categorical reparameterization with gumbel-softmax. In *International Conference on Learning Representations*.

[23] Matthew Johnson, David K Duvenaud, Alex Wiltschko, Ryan P Adams, and Sandeep R Datta. 2016. Composing graphical models with neural networks for structured representations and fast inference. In *Neural Information Processing Systems*.

[24] Eric Jones, Travis Oliphant, Pearu Peterson, et al. 2001. SciPy: Open source scientific tools for Python. (2001). http://www.scipy.org/

[25] Norman P Jouppi, Cliff Young, Nishant Patil, David Patterson, Gaurav Agrawal, Raminder Bajwa, Sarah Bates, Suresh Bhatia, Nan Boden, Al Borchers, Rick Boyle, Pierre-luc Cantin, Clifford Chao, Chris Clark, Jeremy Coriell, Mike Daley, Matt Dau, Jeffrey Dean, Ben Gelb, Tara Vazir Ghaemmaghami, Rajendra Gottipati, William Gulland, Robert Hagmann, C Richard Ho, Doug Hogberg, John Hu, Robert Hundt, Dan Hurt, Julian Ibarz, Aaron Jaffey, Alek Jaworski, Alexander Kaplan, Harshit Khaitan, Andy Koch, Naveen Kumar, Steve Lacy, James Laudon, James Law, Diemthu Le, Chris Leary, Zhuyuan Liu, Kyle Lucke, Alan Lundin, Gordon MacKean, Adriana Maggiore, Maire Mahony, Kieran Miller, Rahul Nagarajan, Ravi Narayanaswami, Ray Ni, Kathy Nix, Thomas Norrie, Mark Omernick, Narayana Penukonda, Andy Phelps, Jonathan Ross, Matt Ross, Amir Salek, Emad Samadiani, Chris Severn, Gregory Sizikov, Matthew Snelham, Jed Souter, Dan Steinberg, Andy Swing, Mercedes Tan, Gregory Thorson, Bo Tian, Horia Toma, Erick Tuttle, Vijay Vasudevan, Richard Walter, Walter Wang, Eric Wilcox, and Doe Hyun Yoon. 2017. In-Datacenter Performance Analysis of a Tensor Processing Unit. *arXiv preprint arXiv:1704.04760* (2017).

[26] Diederik P Kingma, Tim Salimans, Rafal Jozefowicz, Xi Chen, Ilya Sutskever, and Max Welling. 2016. Improving Variational Inference with Inverse Autoregressive Flow. In *Neural Information Processing Systems*.

[27] Diederik P Kingma and Max Welling. 2014. Auto-Encoding Variational Bayes. In *International Conference on Learning Representations*.

[28] Alp Kucukelbir, Dustin Tran, Rajesh Ranganath, Andrew Gelman, and David M Blei. 2017. Automatic Differentiation Variational Inference. *Journal of Machine Learning Research* 18, 14

(2017), 1–45.

[29] Chris J. Maddison, Andriy Mnih, and Yee Whye Teh. 2017. The Concrete Distribution: A Continuous Relaxation of Discrete Random Variables. In *International Conference on Learning Representations*.

[30] George Marsaglia and Wai Wan Tsang. 2000. A simple method for generating gamma variables. *ACM Transactions on Mathematical Software (TOMS)* 26, 3 (2000), 363–372.

[31] Tetsuzo Morimoto. 1963. Markov processes and the H-theorem. *Journal of the Physical Society of Japan* 18, 3 (1963), 328–331.

[32] Alfred Müller. 1997. Integral probability metrics and their generating classes of functions. *Advances in Applied Probability* 29, 2 (1997), 429–443.

[33] Kevin P Murphy. 2012. *Machine Learning: a Probabilistic Perspective.* MIT press.

[34] Christian Naesseth, Francisco Ruiz, Scott Linderman, and David M Blei. 2017. Reparameterization Gradients through Acceptance-Rejection Sampling Algorithms. In *Artificial Intelligence and Statistics*.

[35] Harald Niederreiter. 1978. Quasi-Monte Carlo methods and pseudo-random numbers. *Bull. Amer. Math. Soc.* 84, 6 (1978), 957–1041.

[36] George Papamakarios, Theo Pavlakou, and Iain Murray. 2017. Masked Autoregressive Flow for Density Estimation. In *Neural Information Processing Systems*.

[37] Benjamin Peterson. [n. d.]. six: Python 2 and 3 compatibility utilities. https://github.com/benjaminp/six. ([n. d.]).

[38] Pyro Developers. 2017. Pyro. https://github.com/pyro/pyro. (2017).

[39] Pytorch Developers. 2017. Pytorch. https://github.com/pytorch/pytorch. (2017).

[40] Rajesh Ranganath, Sean Gerrish, and David Blei. 2014. Black box variational inference. In *Artificial Intelligence and Statistics*.

[41] Tim Salimans, Andrej Karpathy, Xi Chen, and Diederik P Kingma. 2017. PixelCNN++: Improving the PixelCNN with Discretized Logistic Mixture Likelihood and Other Modifications. In *International Conference on Learning Representations*.

[42] John Schulman, Nicolas Heess, Theophane Weber, and Pieter Abbeel. 2015. Gradient Estimation Using Stochastic Computation Graphs. In *Neural Information Processing Systems*.

[43] Jiaxin Shi, Jianfei Chen, Jun Zhu, Shengyang Sun, Yucen Luo, Yihong Gu, and Yuhao Zhou. 2017. ZhuSuan: A Library for Bayesian Deep Learning. *arXiv preprint arXiv:1709.05870* (2017).

[44] Michael Spivak. [n. d.]. A Comprehensive Introduction to Differential Geometry, Vol. III. *AMC* 10 ([n. d.]), 12.

[45] Dustin Tran and David Blei. 2017. Implicit Causal Models for Genome-wide Association Studies. *arXiv preprint arXiv:1710.10742* (2017).

[46] Dustin Tran, Matthew D Hoffman, Rif A Saurous, Eugene Brevdo, Kevin Murphy, and David M Blei. 2017. Deep Probabilistic Programming. In *International Conference on Learning Representations*.

[47] Dustin Tran, Alp Kucukelbir, Adji B Dieng, Maja Rudolph, Dawen Liang, and David M Blei. 2016. Edward: A library for probabilistic modeling, inference, and criticism. *arXiv preprint arXiv:1610.09787* (2016).

[48] Aäron van den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex Graves, Nal Kalchbrenner, Andrew Senior, and Koray Kavukcuoglu. 2016. WaveNet: A Generative Model for Raw Audio. *arXiv preprint arXiv:1609.03499* (2016).

[49] Aaron van den Oord, Nal Kalchbrenner, and Koray Kavukcuoglu. 2016. Pixel recurrent neural networks. In *International Conference on Machine Learning*.

[50] Stéfan van der Walt, S Chris Colbert, and Gael Varoquaux. 2011. The NumPy array: a structure for efficient numerical computation. *Computing in Science & Engineering* 13, 2 (2011), 22–30.

[51] Larry Wasserman. 2013. *All of Statistics: A concise Course in Statistical Inference.* Springer Science & Business Media.