

1 Introduction

Deep neural networks have repeatedly surpassed state of the art benchmarks in various use cases, and they have been deployed in many commercial and industry applications. The surge in demand for deep learning has motivated the development of neural network libraries such as tensorflow and pytorch as well as greater abstractions such as keras, which allow more black box model training and development. The associated demand for data scientists and machine learning engineers have led to “transplants” from other fields migrating into machine learning, and much of the knowledge of deep learning has been distilled into the composition of functional, accessible components.

In contrast, the question of what deep learning models are actually doing is unclear and perhaps ill-posed. We seek to present partial answers presented from the fields of statistics, probabilistic graphical models, optimization, and operator theory. Many of these connections are well-known in their respective fields or at least are heavily implicit in their formulations. We are especially interested in the limitations of these interpretations and their intersections, especially where a secondary interpretation offers insight beyond the scope of a more restrictive interpretation.

We begin by describing the exponential family and special properties of the regression of residuals for learning exponential family models, thereby fully describing the output layer of most machine learning models. Second, we describe connect a few very simple probabilistic graphical models with well-known machine learning models, providing an interpretation of pairwise weights for intermediate deep learning layers. Third, we explore the optimization perspective of machine learning, starting with the probabilistic branch of optimization and then generalizing past by exploring other operators which provide an interpretation of the rectified linear unit. Finally, we transition into an operator perspective and the deep equilibrium model perspective of deep learning and tie it into the particle perspective of our layers acting as operators, transforming our spaces.

2 Exponential Families

Exponential families are deeply intertwined into the concept of entropy and energy and can act as a probabilistic building block for density estimation. In fact, both regression and classification are special cases of target distributions defined by a specific exponential family. Moreover, exponential families provide a more specific definition of residuals, and gradient descent of exponential family losses can be shown to directly update the natural parameters. The form for exponential family of distribution of distributions is derived by the Pitman-Koopman-Dermois Theorem. The proof by Koopman built off Fisher’s ideas about sufficient statistics to express a family of distributions which can be defined by the sufficient statistic, assuming certain other restrictions.

$$\begin{aligned} P(x|\eta) &= g(\eta)h(x)\exp\{\langle u(x), \eta \rangle\} \\ P(x|\eta) &= h(x)\exp\{\langle u(x), \eta \rangle - G(\eta)\} \end{aligned} \tag{1}$$

where $u(x)$ is the sufficient statistic (which must be additive) for a single element, η is the natural parameter for that sufficient statistic, $h(x)$ is an innate measure of X which is completely independent of the choice of the natural parameter, $g(\eta)$ is the normalization constant for the Pitman-Koopman-Dermois form, and $G(\eta)$ is an alternative form for the normalization constant which is more consistent with the log factor potential interpretation.

We perform an equivalent derivation for both cases, with $\frac{1}{g(\eta)} = \exp \{G(\eta)\}$:

$$\begin{aligned}
\int P(x|\eta) dx &= \int g(\eta) h(x) \exp \{ \langle u(x), \eta \rangle \} dx \\
1 &= g(\eta) \int h(x) \exp \{ \langle u(x), \eta \rangle \} dx \\
\frac{1}{g(\eta)} &= \int h(x) \exp \{ \langle u(x), \eta \rangle \} dx \\
\int P(x|\eta) dx &= \int h(x) \exp \{ \langle u(x), \eta \rangle - G(\eta) \} dx \\
1 &= \int h(x) \exp \{ \langle u(x), \eta \rangle - G(\eta) \} dx \\
\exp \{G(\eta)\} &= \int h(x) \exp \{ \langle u(x), \eta \rangle \} dx
\end{aligned} \tag{2}$$

++++

One very important property relates the natural parameter to the expected sufficient statistic:

$$\begin{aligned}
\exp \{G(\eta)\} &= \int h(x) \exp \{ \langle u(x), \eta \rangle \} dx \\
\nabla_\eta \exp \{G(\eta)\} &= \frac{\int u(x) h(x) \exp \{ \langle u(x), \eta \rangle \} dx}{\int h(x) \exp \{ \langle u(x), \eta \rangle \} dx} \\
&= \frac{\int u(x) h(x) \exp \{ \langle u(x), \eta \rangle \} dx}{\exp \{G(\eta)\}} < ++ > \\
\nabla_\eta \exp \{G(\eta)\} &= E_{x \sim P(x|\eta)} [u(x)] < ++ > \\
\nabla_\eta [-\log g(\eta)] &= E_{x \sim P(x|\eta)} [u(x)] < ++ >
\end{aligned} \tag{3}$$

++++

It can be proven that this is a one-to-one mapping which in one direction is known as the canonical-link function, and it has many useful properties. One specific case is when we generate samples from an unknown distribution, in which case the η that maximizes the log probability of seeing those samples would be the η that has the same sufficient statistic (assuming that is possible):

$$\begin{aligned}
\nabla_\eta \frac{1}{N} \sum_{n=1}^N \log P(x_n|\eta) &= \frac{1}{N} \sum_{n=1}^N (u(x_n) - \nabla_\eta G(\eta)) < ++ > \\
&= \frac{1}{N} \sum_{n=1}^N (u(x_n) - E_{x \sim P(x|\eta)} [u(x)]) < ++ > \tag{4}
\end{aligned}$$

$$\begin{aligned}
0 &= \frac{1}{N} \sum_{n=1}^N (u(x_n) - E_{x \sim P(x|\eta)} [u(x)]) \\
\sum_{n=1}^N \frac{u(x_n)}{N} &= E_{x \sim P(x|\eta)} [u(x)]
\end{aligned} \tag{5}$$

i++i

The first form shows that gradient descent on the log probabilities, which is the form of both regression and classification loss functions, with respect to the natural parameters is exactly the same as the difference between the sufficient statistic and the expected sufficient statistic. Hence, the concept of a residual for square regression loss is generalized to residuals of sufficient statistics, and the residual can be used as a gradient update to the natural parameters. This is quite remarkable, and when we consider that the samples are taken from some unknown distribution $p_{data}(x)$, we see that maximizing the negative log likelihood is the same as minimizing the cross-entropy, which is the same as minimizing the KL-divergence:

$$\begin{aligned}
KL(p_{data}(x) \| p(x|\eta)) &= H(p_{data}(x)) - \int p_{data}(x) \log p(x|\eta) dx \\
\min_{\eta} KL(p_{data}(x) \| p(x|\eta)) &= \min_{\eta} - \int p_{data}(x) \log p(x|\eta) dx \\
\nabla_{\eta} KL(p_{data}(x) \| p(x|\eta)) &= \nabla_{\eta} E_{x \sim p_{data}(x)} [-\log p(x|\eta)] \\
&= E_{x \sim p_{data}(x)} [u(x_n)] - E_{x \sim p(x|\eta)} [u(x)]
\end{aligned} \tag{6}$$

i++i

This is called an M-projection, and it intertwines exponential families with the concept of entropy. In fact, Boltzmann's derivation of entropy is a special case with a multinoulli distribution. Intuitively, this connection is sensible, since entropy acts upon log probabilities, and the form of an exponential family distribution is in exponential space.

Finally, in the case of conditional models, or discriminative models, we may have some set of non-linear feature mappings for the input, and we may try to model the natural parameters as a linear combination of the features. In the case of a gaussian conditional distribution, this is known as regression, and for classification (multinomial) logistic regression.

$$\begin{aligned}
\eta &= \sum_j w_j \phi_j(x) \\
\nabla_w \log p(x|w) &= \nabla_w \eta \nabla_{\eta} \log p(x|\eta) \\
&= (\phi_1(x) \quad \cdots \quad \phi_M(x)) \nabla_{\eta} \log p(x|\eta)
\end{aligned} \tag{7}$$

i++i

This is also known as a generalized linear model (GLM), with the canonical link functions defining the nonlinearities. Hence, for these regression models, gradient updates for the weights can be done almost directly through the residuals.

There are many more remarkable properties of exponential families, and due to those, it is preferable to try to model distributions as exponential families. In fact, many derivations for regression and classification can be generalized to any exponential family.

3 Probabilistic Graphical Models

i++i

The properties of exponential families and their connection to KL-divergence and entropy motivates the parametrization of probabilistic graphical models (PGMs), specifically the undirected case, Markov networks.

For the case of positive PGMs, that is PGMs for which the probability density/probability mass is never 0, they can be written in the form of a Gibbs distribution:

$$P_{\Phi}(\mathcal{X}) \propto \prod_i \phi_i(C_i) \quad (8)$$

i++i

where $\mathcal{X} = \{X_1, \dots, X_N\}$, i enumerates the set of cliques of the undirected graph, ϕ_i is a real-valued function called the factor potential, and $C_i = \text{Scope}[\phi_i]$.

This result is known as the Hammersley-Clifford theorem, and for positive distributions it connects the concepts of global conditional independence properties with pairwise conditional independence and local, Markov blanket conditional independence.

In the case of positive distributions, the Gibbs distribution can be written in terms of log factor potentials, and one parametrization called log-linear models involves linear weights along with pre-specified log factor potentials:

$$P_{\Phi}(\mathcal{X}) \propto \prod_i e^{-w_i E_i(D_i)} \quad (9)$$

i++i

where $-\log \phi_i = w_i E_i(D_i)$, and E_i are called energy functions, and these unnormalized parametrizations are called energy-based models.

To discern what these distributions are and what energy is, we note that log-linear models are in fact a type of exponential family distribution, with the energy function for each scope being sufficient statistics, and the weights being the natural parameters. This means that once we define the energy functions, the interpretation of the energy functions is the distribution that minimizes KL divergence (maximizes entropy) while matching expected energies (expected sufficient statistics).

In the case of a single clique potential/energy function for the entire model, the associated exponential family model is known as the Boltzmann distribution. In this specific case, we see that the expected energy is the scaled entropy of the model, up to a term involving the normalization constant:

$$\begin{aligned}
P_{E,\eta}(\mathcal{X}) &\propto e^{-E(\mathcal{X})\eta} \\
H(P_{E,\eta}(\mathcal{X})) &= -\eta \int E(\mathcal{X})P_{E,\eta}(\mathcal{X}) + G(\eta) \\
&= -\eta E_{x \sim P_{E,\eta}(\mathcal{X})} [E(\mathcal{X})] + G(\eta) \\
\frac{1}{\eta} H(P_{E,\eta}(\mathcal{X})) &= E_{x \sim P_{E,\eta}(\mathcal{X})} [E(\mathcal{X})] + \frac{G(\eta)}{\eta}
\end{aligned} \tag{10}$$

i++i

In the case of thermodynamics, they use $\frac{1}{T}$ or $\frac{1}{kT}$ instead of η , and they call it the temperature. If we take the expectation of the log of the unnormalized probability distribution, $\tilde{P}_{E,\eta}(\mathcal{X}) = e^{-E(\mathcal{X})\eta}$, we can ignore the normalization factor, and we get the form of Helmholtz free energy, which motivates the energy functional that is optimized in in variational inference and other PGM inference and learning algorithms.

3.1 Machine learning PGMs

i++i

It turns out even very simple Markov network parametrizations can be very good approximations to useful, real world cases. The simplest log factor potentials would involve polynomials of the scope variables. Somewhat independently, it is preferable to have clique potentials with a small number of scope variables due to the exponential complexity of marginalization, with respect to the number of variables. For that reason, let us mostly consider the case of factor potentials involving two variables, an edge, or a single variable, a vertex:

$$\begin{aligned}
P_{\Phi}(\mathcal{X}) &\propto \exp \left\{ -\sum_{i \sim j} w_{i,j} \langle X_i, X_j \rangle - \sum_i u_i X_i \right\} \\
&= e^{-\mathcal{X}^T W \mathcal{X} - U^T \mathcal{X}}
\end{aligned} \tag{11}$$

i++i

where in the equation we have replaced the quadratic coefficients with a single weight matrix, and the linear coefficients with a single vector.

This is known as the Ising model, in the specific case where all the variables are discrete variables of values $+1$ or -1 , and it is the model for electron spins.

In the more general case where variables may take on values of $\{0, 1\}$ or perhaps continuous values, this is called the Boltzmann machine.

When the PGM can be divided into a bipartite graph, this is called a Restricted Boltzmann Machine, and when both sets of variables have the same size, the weight matrix is an $N \times N$ matrix.

i++i Placeholder, I can't justify that it will be equivalent to a Markov network.

For a Boltzmann machine, let's assume the variable $Y_i \in \{0, 1\}$ is a bernoulli variable, X_j are other nodes in \mathcal{X} , and Nb_{Y_i} is the Markov blanket, that is the undirected neighbors of Y_i . If we isolate the terms in the Boltzmann machine to the ones involving Y_i , we end up with

$$\begin{aligned}
P_{\Phi}(Y_i, Y_i^-) &\propto \exp \left\{ \sum_{j \in Nb_{Y_i}} w_{i,j} \langle Y_i, X_j \rangle + u_i Y_i \right\} \\
&= \exp \left\{ Y_i \sum_{j \in Nb_{Y_i}} (w_{i,j} X_j + u_i) \right\}
\end{aligned} \tag{12}$$

⚡⚡⚡

The exponential family form of a Bernoulli variable is

$$P(Y|\eta) = e^{Y\eta - G(\eta)} \tag{13}$$

⚡⚡⚡

so the functional form of our Boltzmann machine motivates us to treat $\eta = w_{i,j}X_j + u_i$ and assume that is the conditional distribution for $Y_i|Nb_{Y_i}$:

$$\begin{aligned}
\eta &= w_{i,j}X_j + u_i \\
P(Y|\eta) &= E_{Y \sim P(Y|\eta)} [Y] = \sigma(Y)
\end{aligned} \tag{14}$$

⚡⚡⚡

Notably, linear combination of other input is functional form of a feed-forward neural network. A caveat is that we have ignored all factors involving X_j , so technically we would have to pass messages through X_j , which would involve passing messages from the rest of the network.

Comparatively, a feed-forward neural network is a Bayesian discriminative model, with an unspecified distribution for the input distribution. We have to make a stronger assumption that the forward weights define a conditional probability as above. The more formal definition we define a conditional random field from a previous layer to the next layer using a Boltzmann machine, with the first layer being the input nodes.

$$\begin{aligned}
H^0 &= P_{data}(X) \\
P(H^{l+1}|H^l) &= \frac{1}{Z(H^l)} RBM(H^l, H^{l+1}) \\
H^L &= Y
\end{aligned} \tag{15}$$

⚡⚡⚡

Since a conditional random field factors out any potentials involving the input, we do not have to worry about the previous layer's messages. We have had to make an additional assumption that we can define the conditional probability of the current layer, H^l , given the previous layer, H^{l-1} , using a CRF, and this means that we have to assume that the edge potentials connecting to the following layer, H^{l+1} do not contribute to the $P(H^l|H^{l-1})$.

One thing to note is that we are using undirected model as a network fragment of a directed model. The directed version of this network fragment is called a conditional Bayesian network, and for the Bernoulli variable case, we can define what is called a logistic conditional probability distribution (logistic CPD) almost exactly as above ⚡⚡⚡.

In fact, this can be generalized to any exponential family, through a generalized linear model with the parent variables, and we can change the target distribution family by changing the nonlinearity and the range of the target variable. There are a few specific assumptions about the structure we are using here worth explicitly stating: 1) The output variables for a layer are conditionally independent of each other, given the input variables, 2) the CPDs have a property called independence of causal influence, which theoretically speaking means that the information from the parents is combined in an associative (and in our case commutative as well) way, and practical speaking it means we can drop out input.

These properties are very strong, though perhaps intuitively desirable assumptions, and they somewhat fit the model for how neurons communicate with each other. Specifically, we can see the strengths of these connections through the complexity of the model - we can parametrize the conditional bayesian network between two layers by $(N^l + 1) \times N^{l+1}$ weights, whereas for unrestricted Bernoulli distributions, we would need to define a separate conditional probability for every conditional outcome, for a complexity of $(2^{N^l} - 1)^{2^{N^{l+1}}} = \mathcal{O}(2^{N^l N^{l+1}})$. Intuitively, when we combine information about the conditioning variable outcomes in this additive way, the symmetries of addition are compressing information and we are reducing the complexity.

Since neural networks with sigmoid non-linearities can be written in this form and transformers are using multinomial logistic CPDs to define an attention variable, we finally have an interpretation of what the layers of a neural network are.

A specific case of the conditional Bayesian network is a conditional linear Gaussian, where we define the conditional probabilities through linear Gaussians:

$$\begin{aligned} P(X^0) &= \mathcal{N}(X^0 | \mu^0, \Sigma^0) \\ P(X^{l+1} | X^l) &= \mathcal{N}(X^{l+1} | W_{l,l+1} X^l + B^{l+1}, \Sigma^{l+1}) \end{aligned} \quad (16)$$

++++

Usually we assume $\Sigma^{l+1} = (\sigma^{l+1})^2 \mathbf{I}$. The joint probability distribution from a linear Gaussian network can be solved analytically through a forward pass of message passing, and this is the same dynamic programming procedure which is used to solve for the optimal solution of model predictive control problems.

So far we have defined generalized linear models using affine functions of the parent variables. This has allowed us to define models by pairwise edge weights, specifically with log linear factors. What would happen if we relax this condition somewhat, so we are using linear combinations of various non-linear functions? It turns out that Gaussian joint probability distributions can always be described by pairwise edge-weights and variable potentials:

$$\begin{aligned} P(\mathcal{X}) &= \mathcal{N}(\mathcal{X} | \mu | \Sigma) \\ &= \exp \left\{ \left\langle (\mathcal{X}, \mathcal{X}^\top), \left(\Sigma^{-1} \mu - \frac{1}{2} \Sigma^{-1} \right) \right\rangle \right\} \\ &= \exp \left\{ -\frac{1}{2} \mathcal{X}^\top \Lambda \mathcal{X} + \mathcal{X}^\top \Lambda \mu \right\} \end{aligned} \quad (17)$$

i++j

This means that we should be able to do a parametrization of the nodes and edges to define the PGM, even if we don't assume a linear-gaussian bayesian factorization. We can also define a conditional probability distribution this way, so given N input variables, X_1, \dots, X_N , and N output variables, Y_1, \dots, Y_N , we can define a conditional random field through a mean and covariance. This is the PGM for a Gaussian Process:

$$P(Y|X) = GP(\mu(x), K(x, x')) \quad (18)$$

i++j

The benefit is that since we know the conditional joint probability distribution, we can condition on observing subsets of Y :

$$P(Y_{unobs}|X, Y_{obs}) = \mathcal{N}\left(\mu_{unobs} + K_{unobs, obs} \Sigma_{obs, obs}^{-1} (Y_{obs} - \mu_{obs}), \Sigma_{unobs, unobs} - \Sigma_{unobs, obs} \Sigma_{obs, obs}^{-1} \Sigma_{obs, unobs}\right) \quad (19)$$

i++j

A caveat is that this is not exactly a PGM, since the edge weights are dependent on the value of the inputs, whereas with a PGM the parametrizations are independent of values. But from a high level, when defining a CPD from a set of N inputs to N outputs, we have machine learning model where we define it through an $N \times N$ matrix representing the edge weights for an RBM, a conditional Bayesian network through a generalized linear model, or an $N \times N$ matrix representing the covariance matrix for a Gaussian Process.

4 Optimization

i++j

In our exponential families discussion, we saw that our maximum likelihood objectives is equivalent to minimizing the KL divergence, and we are minimizing weights that define our parametrization for our probability distribution. For exponential family models, this amounts to minimizing over the weights the expected energy of the model. When we add L2 weight-regularization, it is equivalent to adding an isotropic Gaussian prior on the weights with covariance $\frac{\lambda}{C} \mathbf{I}$, where C is a multiplicative factor of the log likelihood loss that is often left out, and finding the MAP estimate for the weight. In addition, PGM inference algorithms, having broken out of a discriminative modeling objective, involve optimizing a functional known as the Helmholtz free energy functional, which can be considered a relaxation of KL minimization.

In fact, if we have a GLM model with an associated exponential family distribution, our optimization objective is convex, so we are performing a convex optimization problem. However, the deeper layers change it into a non-convex problem.

There are some interesting tools in convex optimization that may be of help here.

One is under certain restrictions, we can compose functions together into convex functions. We can compose convex/concave functions with real-valued monotonic convex/concave into convex or concave functions, where the result depends on whether the scalar functions is nonincreasing or nondecreasing. We

can also do composition with vector-value component-wise monotonic functions and and guarantee convexity/concavity. This is suggestively useful, since affine functions are both linear and concave, and neural networks are built with compositions of affine functions and non-linearities - though the sigmoid nonlinearity is neither convex nor concave.

Another way is that we can relax the definition of convexity to quasiconvexity, which guarantees weakened properties of convexity. Essentially this allows a function to become concave, but only if it is decreasing, although we have to be careful about what decreasing means in multi-dimensional spaces. Practically speaking, quasiconvexity is a property about the sublevel sets being convex. Moreover, many optimization algorithm proofs assume that the sublevel sets are bounded, which means that whatever loss value we are currently at, we only have to search over a bounded, convex set to find the minimum. In addition, composition of functions does apply to quasiconvex functions as well. This is suggestive of the ability of optimization to provide an interpretation of multi-layered networks.

Though we are able to convert PGMs into optimization objectives, we quickly find departures from this ideal in machine learning. One example is the Rectified linear unit (ReLU), which on inspection does not seem to have a well-known exponential family distribution analog. A second is that often additional regularization terms are added, such as smooth differentiability regularization or consistency regularization, that do not known probability distributions.

In contrast, regularization is a very well-developed theory in optimization. The primary optimization objectives and each of the regularization terms become what is called a multicriterion optimization. Constrained optimization is one case of multicriterion optimization, where bounds on the constraints function values define the additional criteria, and for Lagrangian optimization, the regularization constants have to be optimized as well, through the dual problem. Regularization, or softly constrained problems, pre-define the regularization constants.

Overall, the field of optimization has developed algorithms for optimizing the dual problem to turn a hard-constrained problem into a soft-constrained optimization problem, as well as primal-dual algorithms.

The interpretation of ReLU functions is somewhat more complicated. As discussed above, we might have some hard-constraint on the space where we are optimizing over. Normally, we would use gradient descent, but once we reach the constraint boundaries, our objective function has changed into something that is not differentiable. One method for avoiding this non-differentiability is to project the update back into the constrained set of legal weight values, and these two steps are what is called the projected gradient method. In fact, vanilla gradient descent with stepsize t_k is equivalent to minimizing a quadratic model:

$$\theta^{k+1} := \arg \min_{\theta} f(\theta^k) + \langle \nabla_{\theta} f(\theta^k), \theta - \theta^k \rangle + \frac{t_k}{2} \|\theta - \theta^k\|^2 = \theta^k - t_k \nabla_{\theta} f(\theta^k) \quad (20)$$

i++i

If we minimize that quadratic model over a set C , we define our projection operator:

$$\begin{aligned}
P_C(\theta) &= \arg \min_{\theta} \delta(\theta \geq 0) \\
P_C(\theta^k - t_k \nabla_{\theta} f(\theta^k))
\end{aligned} \tag{21}$$

i++j

One interesting projection we can apply is projection to the non-negative real numbers - that is in fact the ReLu function. There are useful optimization implications for this function - when dealing with Lagrange multipliers, for our inequality constraints we have to constrain our lagrange multipliers to have non-negative values.

Just as hard constraints are the analog to regularization, we can soften the hard constraints of a projection operator to regularization functions, like L1 regularization or other, possibly non-differentiable functions. This is called the proximal operator:

$$\begin{aligned}
prox_g(x) &= \arg \min_u g(u) + \frac{1}{2} \|x - u\|^2 \\
\theta^{k+1} &:= prox_{t_k g}(\theta^k - t_k \nabla_{\theta} f(\theta^k))
\end{aligned} \tag{22}$$

i++j

The use of these operators as a secondary step on top of a gradient descent step opens up a more powerful field of optimization theory and algorithms, including the fast iterative shrinking and thresholding algorithm (FISTA), of which Nesterov momentum is a special case.

From an interpretative standpoint, we can see that components in our model that involve projections and other operators may be perfectly valid as an alternative way to perform an optimization step without having to perform differentiation.

5 Operator Theory

i++j