# A Geometric View of Conjugate Priors

**Arvind Agarwal**
Department of Computer Science
University of Maryland
College Park, Maryland USA
arvinda@cs.umd.edu

**Hal Daumé III**
Department of Computer Science
University of Maryland
College Park, Maryland USA
hal@cs.umd.edu

## Abstract

In Bayesian machine learning, conjugate priors are popular, mostly due to mathematical convenience. In this paper, we show that there are deeper reasons for choosing a conjugate prior. Specifically, we formulate the conjugate prior in the form of Bregman divergence and show that it is the inherent geometry of conjugate priors that makes them appropriate and intuitive. This geometric interpretation allows one to view the hyperparameters of conjugate priors as the *effective* sample points, thus providing additional intuition. We use this geometric understanding of conjugate priors to derive the hyperparameters and expression of the prior used to couple the generative and discriminative components of a hybrid model for semi-supervised learning.

## 1 Introduction

In probabilistic modeling, a practitioner typically chooses a likelihood function (model) based on her knowledge of the problem domain. With limited training data, a simple maximum likelihood estimation (MLE) of the parameters of this model will lead to overfitting and poor generalization. One can regularize the model by adding a prior, but the fundamental question is: which prior? We give a turn-key answer to this problem by analyzing the underlying *geometry* of the likelihood model, and suggest choosing the unique prior with the same geometry as the likelihood. This unique prior turns out to be the *conjugate* prior, in the case of the exponential family. This provides justification beyond "computational convenience" for using the conjugate prior in machine learning and data mining applications.

In this work, we give a geometric understanding of the maximum likelihood estimation method and a geometric argument in the favor of using conjugate priors. Empirical evidence showing the effectiveness of the conjugate priors can be found in our earlier work [1]. In Section 4.1, first we formulate the MLE problem into a completely geometric problem with no explicit mention of probability distributions. We then show that this geometric problem carries a geometry that is inherent to the structure of the likelihood model. For reasons given in Sections 4.3 and 4.4, when considering the prior, it is important that one uses the same geometry as likelihood.

Using the same geometry also gives the closed-form solution for the maximum-a-posteriori (MAP) problem. We then analyze the prior using concepts borrowed from the information geometry. We show that this geometry induces the *Fisher information metric* and *1-connection*, which are respectively, the natural metric and connection for the exponential family (Section 5). One important outcome of this analysis is that it allows us to treat the hyperparameters of the conjugate prior as the effective sample points drawn from the distribution under consideration. We finally extend this geometric interpretation of conjugate priors to analyze the hybrid model given by [7] in a purely geometric setting, and justify the argument presented in [1] (i.e. a *coupling prior* should be conjugate) using a much simpler analysis (Section 6). Our analysis couples the discriminative and generative components of hybrid model using the Bregman divergence which reduces to the coupling prior given in [1]. This analysis avoids the *explicit* derivation of the hyperparameters, rather automatically gives the hyperparameters of the conjugate prior along with the expression.

## 2 Motivation

Our analysis is driven by the desire to understand the geometry of the conjugate priors for the exponential families. We motivate our analysis by asking ourselves the following question: Given a parametric model $p(x; \theta)$ for the data likelihood, and a prior on its parameters $\theta$, $p(\theta; \alpha, \beta)$; what should the hyperparameters $\alpha$ and $\beta$ of the prior encode? We know that $\theta$ in the likelihood model is the estimation of the parameter using the given data points. In other words, the estimated parameter fits the model according to the given data while the prior on the parameter provides the generalization. This generalization is enforced by some prior belief encoded in the hyperparameters. Unfortunately, one does not know what is the likely value of the parameters; rather one might have some belief in what *data points* are likely to be sampled from the model. Now the question is: Do the hyperparameters encode this belief in the parameters in terms of the sampling points? Our analysis shows that the hyperparameters of the conjugate prior is nothing but the effective sampling points. In case of non-conjugate priors, the interpretation of hyperparameters is not clear.

A second motivation is the following geometric analysis. Before we go into the problem, consider two points in the
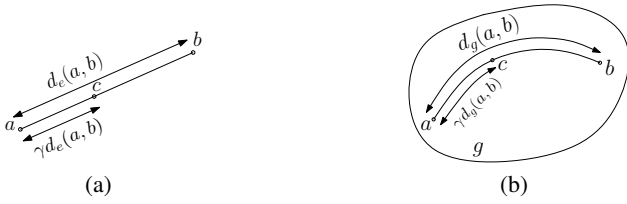
Figure 1: Interpolation of two points $a$ and $b$ using (a) Euclidean geometry, and (b) non-Euclidean geometry. Here geometry is defined by the respective distance/divergence functions $d_e$ and $d_g$. It is important to notice that the divergence is a generalized notion of the distance in the non-Euclidean spaces, in particular, in the spaces of the exponential family statistical manifolds. In these spaces, it is the divergence function that define the geometry.

*Euclidean* space which one would like to interpolate using a parameter $\gamma \in [0, 1]$. A natural way to do so is to interpolate them linearly i.e., connect two points using a straight line, and then find the interpolating point at the desired $\gamma$, as shown in Figure 1(a). This interpolation scheme does not change if we move to a non-Euclidean space. In other words, if we were to interpolate two points in a non-Euclidean space, we would find the interpolating point by connecting two points by a geodesic (an equivalent to the straight line in the non-Euclidean space) and then finding the point at the desired $\gamma$, shown in Figure 1(b).

This situation arises when one has two models, and wants to build a better model by interpolating them. This exact situation is encountered in [7] where the objective is to build a hybrid model by interpolating (or coupling) discriminative and generative models. Agarwal et.al. [1] couples these two models using the conjugate prior, and empirically shows using a conjugate prior for the coupling outperforms the original choice [7] of a Gaussian prior. In this work, we find the hybrid model by interpolating the two models using the *inherent geometry*[1] of the space (interpolate along the geodesic in the space defined by the inherent geometry) which automatically results in the conjugate prior along with its hyperparameters. Our analysis and the analysis of Agarwal et al. lead to the same result, but ours is much simpler and naturally extends to the cases where one wants to couple more than two models. One big advantage of our analysis is that unlike prior approaches [1], we need not know the expression and the hyperparameters of the prior in advance. They are automatically derived by the analysis. Our analysis only requires the inherent geometry of the models under consideration and the interpolation parameters. No explicit expression of the coupling prior is needed.

## 3 Exponential Family

In this section, we review the exponential family. The exponential family is a set of distributions, whose probability



Figure 2: Duality between mean parameters and natural parameters.

density function can be expressed in the following form:

$$p(x; \theta) = p_o(x)\exp(\langle \theta, \phi(x) \rangle - G(\theta)) \qquad (1)$$

here $\phi(x) : \mathcal{X}^m \rightarrow \mathbb{R}^d$ is a vector *potentials* or *sufficient statistics* and $G(\theta)$ is a normalization constant or *log-partition function*. With the potential functions $\phi(x)$ fixed, every $\theta$ induces a particular member $p(x; \theta)$ of the family. In our framework, we deal with exponential families that are *regular* and have the *minimal representation*[9].

One important property of the exponential family is the existence of conjugate priors. Given any member of the exponential family in (1), the *conjugate prior* is a distribution over its *parameters* with the following form:

$$p(\theta | \alpha, \beta) = m(\alpha, \beta) \exp(\langle \theta, \alpha \rangle - \beta G(\theta)) \qquad (2)$$

here $\alpha$ and $\beta$ are hyperparameters of the conjugate prior. Importantly, the function $G(\cdot)$ is the same between the exponential family member and its conjugate prior.

A second important property of exponential family member is that log-partition function $G$ is convex and defined over the convex set $\Theta := \{\theta \in \mathbb{R}^d : G(\theta) < \infty\}$; and since it is convex over this set, it induces a Bregman divergence [3][2] on the space $\Theta$.

Another important property of the exponential family is the *one-to-one* mapping between the *canonical parameters* $\theta$ and the so-called "*mean parameters*" which we denote by $\mu$. For each canonical parameter $\theta \in \Theta$, there exists a mean parameter $\mu$, which belongs to the space $\mathcal{M}$ defined as:

$$\mathcal{M} := \left\{ \mu \in \mathbb{R}^d : \mu = \int \phi(x)p(x; \theta) \, dx \quad \forall \theta \in \Theta \right\} \qquad (3)$$

It is easy to see that $\Theta$ and $\mathcal{M}$ are dual spaces, in the sense of Legendre (conjugate) duality because of the following relationship between the log-partition function $G(\theta)$ and the expected value of the sufficient statistics $\phi(x)$: $\nabla G(\theta) = \mathbb{E}(\phi(x)) = \mu$. In Legendre duality, we know that two spaces $\Theta$ and $\mathcal{M}$ are dual of each other if for each $\theta \in \Theta$, $\nabla G(\theta) = \mu \in \mathcal{M}$. We call the function in the dual space $\mathcal{M}$ to be $F$ i.e., $F = G^*$. A pictorial representation of the duality between canonical parameter space $\Theta$ and mean parameter space $\mathcal{M}$ is given in Figure 2.

---

[1]In exponential family statistical manifold, inherent geometry is defined by the divergence function because it is the divergence function that induces the metric structure and connection of the manifold. Refer [2] for more details.
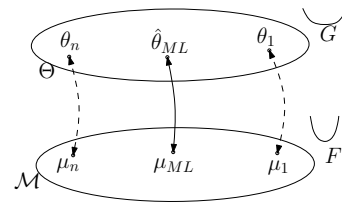
[2]Two important points to note about Bregman divergence are: 1) For dual spaces $F$ and $G$, $B_F(p\|q) = B_G(q^*\|p^*)$, where $p^*$ and $q^*$ are the conjugate duals of $p$ and $q$ respectively. 2) Bregman divergence is not symmetric i.e., in general, $B_F(p\|q) \neq B_F(q\|p)$, therefore it is important what directions these divergences are measured in.

In our analysis, we will need the Bregman divergence over $\phi(x)$ which can be obtained by showing that an augmented $\mathcal{M}$ contains all possible $\phi(x)$. In order to define the Bregman divergence over all $\phi(x)$, we define a new set of mean parameters w.r.t. all probability distributions (*not* only w.r.t. exponential family distributions): $\mathcal{M}^+ := \{\mu \in \mathbb{R}^d : \mu = \int \phi(x)p(x)\,dx \quad \text{s.t.} \int p(x)\,dx = 1\}$.

Note that $\mathcal{M}^+$ is the convex hull of $\phi(x)$ thus contains all $\phi(x)$. We know from (see Theorem 3.3, [10]) that $\mathcal{M}$ is the *interior* of $\mathcal{M}^+$. Now we augment $\mathcal{M}$ with the boundary of $\mathcal{M}^+$ and $\Theta$ with the canonical parameters (limiting distributions) that will generate the mean parameters corresponding to this boundary. We know (see Theorem 2, [9]) that such parameters exist. Call these new sets $\mathcal{M}^+$ and $\Theta^+$ respectively. We also know [9] that $\Theta^+$ and $\mathcal{M}^+$ are conjugate dual of each other (for boundary, duality exists in the limiting sense) i.e., Bregman divergence is defined over the entire $\mathcal{M}^+$.

In the following discussion, $\mathcal{M}$ and $\Theta$ will denote the closed sets i.e. $\mathcal{M}^+$ and $\Theta^+$ respectively.

## 4 Likelihood, Prior and Geometry

In this section, we first formulate the ML problem into a Bregman median problem (Section 4.1) and then show that corresponding MAP (maximum-a-posteriori) problem can also be converted into a Bregman median problem (Section 4.3). The MAP Bregman median problem consists of two parts: a likelihood model and a prior. We argue (Sections 4.3 and 4.4) that a Bregman median problem makes sense only when both of these parts have the same geometry. Having the same geometry amounts to having the same log-partition function leading to the property of conjugate priors.

### 4.1 Likelihood in the form of Bregman Divergence

Following [5], we can write the distributions belonging to the exponential family (1) in terms of Bregman divergence [3]:

$$\log p(x;\theta) = \log p_o(x) + F(x) - B_F(x\|\nabla G(\theta)) \quad (4)$$

This representation of likelihood in the form of Bregman divergence gives insight in the geometry of the likelihood function. Gaining the insight into the exponential family distributions, and establishing a meaningful relationship between likelihood and prior is the primary objective of this work.

In learning problems, one is interested in estimating the parameters $\theta$ of the model which results in low generalization error. Perhaps the most standard estimation method is *maximum likelihood* (ML). The ML estimate, $\hat{\theta}_{ML}$, of a set of $n$ i.i.d. training data points $\mathcal{X} = \{x_1, \ldots x_n\}$ drawn from the exponential family is obtained by solving the following problem: $\hat{\theta}_{ML} = \max_{\theta \in \Theta} \log p(\mathcal{X};\theta) = \max_{\theta \in \Theta} \sum_{i=1}^{n} \log p(x_i;\theta)$.

**Theorem 1.** *Let $\mathcal{X}$ be a set of $n$ i.i.d. training data points drawn from the exponential family distribution with the log partition function $G$, $F$ be the dual function of $G$, then dual of ML estimate ($\hat{\theta}_{ML}$) of $\mathcal{X}$ under the assumed exponential family model solves the following Bregman median problem:* $\hat{\mu}_{ML} = \min_{\mu \in \mathcal{M}} \sum_{i=1}^{n} B_F(x_i\|\mu)$.

[3]For the simplicity of the notations we will use $x$ instead of $\phi(x)$ assuming that $x \in \mathbb{R}^d$. This does not change the analysis

*Proof.* Proof is straightforward. Using (4) in MLE problem $\max_{\theta \in \Theta} \sum_{i=1}^{n} \log p(x_i;\theta)$, and ignoring terms that do not depend on $\theta$:

$$\hat{\theta}_{ML} = \min_{\theta \in \Theta} \sum_{i=1}^{n} B_F(x_i\|\nabla G(\theta)) \quad (5)$$

which using the expression $\nabla G(\theta) = \mu$ gives the desired result. $\qquad\square$

The above theorem converts the problem of maximizing the log likelihood $\log p(\mathcal{X};\theta)$ into an equivalent problem of minimizing the corresponding Bregman divergences which is nothing but a *Bregman median* problem, the solution to which is given by $\hat{\mu}_{ML} = \sum_{i=1}^{n} x_i$. ML estimate $\hat{\theta}_{ML}$ can now be computed using the expression $\nabla G(\theta) = \mu$, $\hat{\theta}_{ML} = (\nabla G)^{-1}(\hat{\mu}_{ML})$.

**Lemma 1.** *If $x$ is the sufficient statistics of the exponential family with the log partition function $G$, and $F$ is the dual function of $G$ defined over the mean parameter space $\mathcal{M}$ then (1) $x \in \mathcal{M}$; (2) there exists a $\theta \in \Theta$, such that $x^* = \theta$.*

*Proof.* (1) By construction of $\mathcal{M}$, we know $x \in \mathcal{M}$. (2) From duality of $\mathcal{M}$ and $\Theta$, for every $\mu \in \mathcal{M}$, there exists a $\theta \in \Theta$ such that $\theta = \mu^*$, and since $x \in \mathcal{M}$, which implies $x^* = \theta$. $\qquad\square$

**Corollary 1 (ML as Bregman Median).** *Let $G$ and $\mathcal{X}$ be defined as earlier, $\theta_i$ be the dual of $x_i$, then ML estimation, $\hat{\theta}_{ML}$ of $\mathcal{X}$ solves the following optimization problem:*

$$\hat{\theta}_{ML} = \min_{\theta \in \Theta} \sum_{i=1}^{n} B_G(\theta\|\theta_i) \quad (6)$$

*Proof.* Proof directly follows from Lemma 1 and Theorem 1. From Lemma 1, we know that $x_i^* = \theta_i$. Using Theorem 1 and expression $B_F(x_i\|\mu) = B_G(\theta\|x_i^*) = B_G(\theta\|\theta_i)$ gives the desired result. $\qquad\square$

The above expression requires us to find a $\theta$ so that divergence from $\theta$ to other $\theta_i$ is minimized. Now note that $G$ is what defines this divergence and hence the geometry of the $\Theta$ space (as discussed earlier in Section 2). since $G$ is the log partition function of an exponential family, **it is the log-partition function that determines the geometry of the space**. We emphasize that divergence is measured from the parameter being estimated to other parameters $\theta_i$(s), as shown in Figure 3.

### 4.2 Conjugate Prior in the form of Bregman Divergence

We now give an expression similar to the likelihood for the conjugate prior (ignoring the term $\log m(\alpha, \beta)$):

$$\log p(\theta|\alpha, \beta) = \beta(\langle \theta, \tfrac{\alpha}{\beta} \rangle - G(\theta)) \quad (7)$$

which can be written in the form of Bregman divergence by a direct comparison to (1), replacing $x$ with $\alpha/\beta$.

$$\log p(\theta|\alpha, \beta) = \beta \left( F\left(\frac{\alpha}{\beta}\right) - B_F\left(\frac{\alpha}{\beta}\|\nabla G(\theta)\right) \right) \quad (8)$$

The expression for the joint probability of data and parameters (combining all terms that do not depend on $\theta$ in const) is given by:

$$\log p(x, \theta | \alpha, \beta) = \text{const} - B_F(x \| \mu) - \beta B_F\left(\frac{\alpha}{\beta} \| \mu\right) \quad (9)$$

### 4.3 Geometric Interpretation of Conjugate Prior

In this section we give a geometric interpretation of the term $B_F(x \| \mu) + \beta B_F\left(\frac{\alpha}{\beta} \| \mu\right)$ from (9).

**Theorem 2 (MAP as Bregman median).** *Given a set $\mathcal{X}$ of $n$ i.i.d examples drawn from the exponential family distribution with the log partition function $G$ and a conjugate prior as in (8), MAP estimation of parameters is $\hat{\theta}_{MAP} = \hat{\mu}^*_{MAP}$ where $\hat{\mu}_{MAP}$ solves the following problem:*

$$\hat{\mu}_{MAP} = \min_{\mu \in \mathcal{M}} \sum_{i=1}^{n} B_F(x_i \| \mu) + \beta B_F\left(\frac{\alpha}{\beta} \| \mu\right) \quad (10)$$

*which admits the following solution:* $\hat{\mu}_{MAP} = \frac{\sum_{i=1}^{n} x_i + \alpha}{n + \beta}$.

*Proof.* Proof is a direct result of applying the definition of Bregman divergence on (9) for all $n$ points. $\square$

The above solution gives a natural interpretation of MAP estimation. One can think of prior as $\beta$ number of extra points at position $\alpha/\beta$. $\beta$ works as the effective sample size of the prior which is clear from the following expression of the dual of the $\hat{\theta}_{MAP}$:

$$\hat{\mu}_{MAP} = \frac{\sum_{i=1}^{n} x_i + \sum_{i=1}^{\beta} \frac{\alpha}{\beta}}{n + \beta} \quad (11)$$

The expression (10) is analogous to (5) in the sense that both are defined in the dual space $\mathcal{M}$. One can convert (10) into an expression similar to (6) in the dual space which is again a Bregman median problem in the parameter space.

$$\hat{\theta}_{MAP} = \min_{\theta \in \Theta} \sum_{i=1}^{n} B_G(\theta \| \theta_i) + \sum_{i=1}^{\beta} B_G\left(\theta \| \left(\frac{\alpha}{\beta}\right)^*\right) \quad (12)$$

here $\left(\frac{\alpha}{\beta}\right)^* \in \Theta$ is dual of $\frac{\alpha}{\beta}$. The above problem is a Bregman median problem of $n + \beta$ points, $\{\theta_1, \ldots \theta_n, \underbrace{(\alpha/\beta)^*, \ldots, (\alpha/\beta)^*}_{\beta \text{ times}}\}$, as shown in Figure 3 (left).

A geometric interpretation is also shown in Figure 3. When the prior is conjugate to the likelihood, they both have the same log-partition function (Figure 3, left). Therefore they induce the same Bregman divergence. Having the same divergence means that distances from $\theta$ to $\theta_i$ (in likelihood) and the distances from $\theta$ to $(\alpha/\beta)^*$ are measured with the same divergence function, yielding the same geometry for both spaces.

It is easier to see using the median formulation of the MAP estimation problem that one must choose a prior that is conjugate. If one chooses a conjugate prior, then the distances among all points are measured using the same function. It is also clear from (11) that in the conjugate prior case, the point induced by the conjugate prior behaves as a sample point $(\alpha/\beta)^*$. A median problem over a space that have different geometries is an ill-formed problem, as discussed further in the next section.
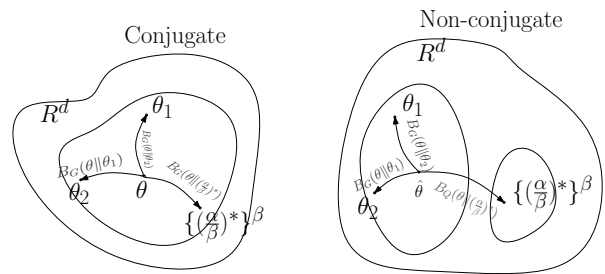


Figure 3: Prior in the conjugate case has the same geometry as the likelihood while in the non-conjugate case, they have different geometries.

### 4.4 Geometric Interpretation of Non-conjugate Prior

We derived expression (12) because we considered the prior conjugate to the likelihood function. Had we chosen a non-conjugate prior with log-partition function $Q$, we would have obtained:

$$\hat{\theta}_{ML} = \min_{\theta \in \Theta} \sum_{i=1}^{n} B_G(\theta \| \theta_i) + \sum_{i=1}^{\beta} B_Q\left(\theta \| \left(\frac{\alpha}{\beta}\right)^*\right). \quad (13)$$

Here $G$ and $Q$ are different functions defined over $\Theta$. Since these are the functions that define the geometry of the space parameter, having $G \neq Q$ is equivalent to consider them as being defined over different (metric) spaces. Here, it should be noted that distance between the sample point $(\theta_i)$ and the parameter $\theta$ is measured using the Bregman divergence $B_G$. On the other hand, the distance between the point induced by the prior $(\alpha/\beta)^*$ and $\theta$ is measured using the divergence function $B_Q$. This means that $(\alpha/\beta)^*$ can *not* be treated as one of the sample points. This tells us that, unlike the conjugate case, belief in the non-conjugate prior can not be encoded in the form of the sample points.

Another problem with considering a non-conjugate prior is that the dual space of $\Theta$ under different functions would be different. Thus, one will not be able to find the alternate expression for (13) equivalent to (10), and therefore not be able to find the closed-form expression similar to (11). This tells us why non-conjugate does not give us a closed form solution for $\hat{\theta}_{MAP}$. A pictorial representation of this is also shown in Figure 3. Note that, unlike the conjugate case, in the non-conjugate case, the data likelihood and the prior both belong to different spaces. We emphasize that it is possible to find the solution of (13) that is, in practice, there is nothing that prohibits the use of non-conjugate prior, however, using the conjugate prior is intuitive, and allows one to treat the hyper-parameters as pseudo data points.

## 5 Information Geometric View

In this section, we show the appropriateness of the conjugate prior from the information geometric angle. In information geometry, $\Theta$ is a statistical manifold such that each $\theta \in \Theta$ defines a probability distribution. This statistical manifold has an inherent geometry, given by a *metric* and an *affine connection*. One natural metric is the Fisher information metric because of its many attractive properties: it is Riemannian and

is invariant under reparameterization (for more details refer [2]).

In exponential family distributions, the Fisher metric $M(\theta)$ is induced by the KL-divergence $KL(\cdot\|\theta)$, which is equivalent to the Bregman divergence defined by the log-partition function. Thus, it is the log-partition function $G$ that induces the Fisher metric, and therefore determines the *natural* geometry of the space. It justifies our earlier argument of choosing the log-partition function to define the geometry. Now if we were to treat the prior as a point on the statistical manifold defined by the likelihood model, the Fisher information metric on the point given by the prior must be same as the one defined on likelihood manifold. This means that the prior must have the same log-partition function as the likelihood i.e., it must be conjugate.

# 6 Hybrid model

In this section, we show an application of our analysis to a common supervised and semi-supervised learning framework. In particular, we consider a generative/discriminative hybrid model [1; 6; 7] that has been shown to be successful in many application. The hybrid model is a mixture of discriminative and generative models, each of which has its own separate set of parameters. These two sets of parameters (hence two models) are combined using a prior called the *coupling prior*. Let $p(y|\mathbf{x}, \theta_d)$ be the discriminative component, $p(\mathbf{x}, y|\theta_g)$ be the generative component and $p(\theta_d, \theta_g)$ be the coupling prior, the joint likelihood of the data and parameters can be written as (combining all three):

$$p(\mathbf{x}, y, \theta_d, \theta_g) = p(\theta_g, \theta_d)p(y|\mathbf{x}, \theta_d)\sum_{y'} p(\mathbf{x}, y'|\theta_g) \quad (14)$$

The most important aspect of this model is the *coupling prior* $p(\theta_g, \theta_d)$, which *interpolates* the hybrid model between two extremes: fully generative when the prior forces $\theta_d = \theta_g$, and fully discriminative when the prior renders $\theta_d$ and $\theta_g$ independent. In non-extreme cases, the goal of the coupling prior is to encourage the generative model and the discriminative model to have similar parameters. It is easy to see that this effect can be induced by many functions. One obvious way is to *linearly* interpolate them as done by [7; 6] using a Gaussian prior (or the Euclidean distance) of the following form:

$$p(\theta_g, \theta_d) \propto \exp\left(-\lambda \|\theta_g - \theta_d\|^2\right) \quad (15)$$

where, when $\lambda = 0$, model is purely discriminative while for $\lambda = \infty$, model is purely generative. Thus $\lambda$ in the above expression is the interpolating parameter, and is same as the $\gamma$ in Section 2. Note that the log of the prior is nothing but the squared Euclidean distance between two sets of parameters.

It has been noted multiple times [4; 1] that a Gaussian prior is not always appropriate, and the prior should instead be chosen according to models being considered. Agarwal et al. [1] suggested using a prior that is conjugate to the generative model. Their main argument for choosing the conjugate prior came from the fact that this provides a closed form solution for the generative parameters and therefore is mathematically convenient. We will show that it is more than convenience that makes conjugate prior appropriate. Moreover, our analysis does not assume anything about the expression and the hyperparameters of the prior beforehand, rather derive them automatically.
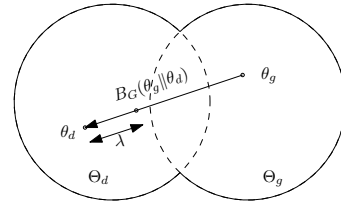


Figure 4: Parameters $\theta_d$ and $\theta_g$ are interpolated using the Bregman divergence

## 6.1 Generalized Hybrid Model

In order to see the effect of the geometry, we present the discriminative and generative models associated with the hybrid model in the Bregman divergence form and obtain their geometry. Following the expression used in [1], the generative model can be written as:

$$p(\mathbf{x}, y|\theta_g) = h(\mathbf{x}, y)\exp(\langle\theta_g, T(\mathbf{x}, y)\rangle - G(\theta_g)) \quad (16)$$

where $T(\cdot)$ is the potential function similar to $\phi$ in (1), now only defined on $(\mathbf{x}, y)$. Let $G^*$ be the dual function of $G$; the corresponding Bregman divergence is given by $B_{G^*}((\mathbf{x}, y)\|\nabla G(\theta_g))$. Solving the generative model independently reduces to choosing a $\theta_g$ from the space of all generative parameters $\Theta_g$ which has a geometry defined by the log-partition function $G$. Similarly to the generative model, the exponential form of the discriminative model is given as:

$$p(y|\mathbf{x}, \theta_d) = \exp(\langle\theta_d, T(\mathbf{x}, y)\rangle - M(\theta_d, \mathbf{x})) \quad (17)$$

Importantly, the sufficient statistics $T$ are the *same* in the generative and discriminative models; such generative/discriminative pairs occur naturally: logistic regression/naive Bayes and hidden Markov models/conditional random fields are examples. However, observe that in the discriminative case, the log partition function $M$ depends on both $\mathbf{x}$ and $\theta_d$ which makes the analysis of the discriminative model (and hence of hybrid model) harder.

## 6.2 Geometry of the Hybrid Model

We simplify the analysis of the hybrid model by rewriting the discriminative model in a a form that makes its underlying geometry obvious. Note that the only difference between the two models is that discriminative model models the conditional distribution while the generative model models the joint distribution. We can use this observation to write the discriminative model in the following alternate form using the expression $p(y|x, \theta) = \frac{p(y, x|\theta)}{\sum_{y'} p(y'x|\theta)}$ and (16):

$$p(y|x, \theta_d) = \frac{h(\mathbf{x}, y)\exp(\langle\theta_d, T(\mathbf{x}, y)\rangle - G(\theta_d))}{\sum_{y'} h(\mathbf{x}, y')\exp(\langle\theta_d, T(\mathbf{x}, y')\rangle - G(\theta_d))} \quad (18)$$

Denote the space of parameters of the discriminative model by $\Theta_d$. It is easy to see that geometry of $\Theta_d$ is defined by $G$ since function $G$ is defined over $\theta_d$. This is same as the geometry of the parameter space of the generative model $\Theta_g$. Now let us define a new space $\Theta_H$ which is the *affine* combination of $\Theta_d$ and $\Theta_g$. Now, $\Theta_H$ will have the same geometry as $\Theta_d$ and $\Theta_g$ i.e., geometry defined by $G$. Now the goal of the hybrid model is to find a $\theta \in \Theta_H$ that maximizes the likelihood of the data under the hybrid model. These two spaces are shown pictorially in Figure 4.

## 6.3 Prior Selection

As mentioned earlier, the coupling prior is the most important part of the hybrid model, which controls the amount of coupling between the generative and discriminative models. There are many ways to do this, one of which is given by [7; 6]. By their choice of Gaussian prior as coupling prior, they implicitly couple the discriminative and generative parameters by the squared Euclidean distance. We suggest coupling these two models by a general prior, of which the Gaussian prior is a special case.

**Bregman Divergence and Coupling Prior:**

Let a general coupling be given by $B_S(\theta_g \| \theta_d)$. Notice the direction of the divergence. We have chosen this direction because the prior is induced on the generative parameters, and it is clear from (12) that parameters on which prior is induced, are placed in the first argument in the divergence function. The direction of the divergence is also shown in Figure 4.

Now we rewrite (8) replacing $\nabla G(\theta)$ by $\theta^*$:

$$\log p(\theta_g | \alpha, \beta) = \beta(F(\frac{\alpha}{\beta}) - B_F(\frac{\alpha}{\beta} \| \theta_g^*)) \qquad (19)$$

Now taking the $\alpha = \lambda\theta_d^*$ and $\beta = \lambda$, we get:

$$p(\theta_g | \lambda\theta_d^*, \lambda) = \exp(\lambda(F(\theta_d^*))) \exp(-\lambda B_F(\theta_d^* \| \theta_g^*)) \qquad (20)$$

For the general coupling divergence function $B_S(\theta_g \| \theta_d)$, the corresponding coupling prior is given by:

$$\exp(-\lambda B_{S^*}(\theta_d^* \| \theta_g^*)) = \exp(-\lambda(F(\theta_d^*)))\, p(\theta_g | \lambda\theta_d^*, \lambda) \qquad (21)$$

The above relationship between the divergence function (left side of the expression) and coupling prior (right side of the expression) allows one to define a Bregman divergence for a given coupling prior and vise versa.

**Coupling Prior for the Hybrid Model:**

We now use (21) to derive the expression for the coupling prior using the geometry of the hybrid model which is given by the log partition function $G$ of the generative model. This argument suggests to couple the hybrid model by the divergence $B_G(\theta_g \| \theta_d)$ which gives the coupling prior as:

$$\exp(-\lambda B_G(\theta_g \| \theta_d)) = p(\theta_g | \lambda\theta_d^*, \lambda) \exp(-\lambda F(\theta_d^*)) \qquad (22)$$

where $\lambda = [0, \infty]$ is the interpolation parameter, interpolating between the discriminative and generative extremes. In dual form, the above expression can be written as:

$$\exp(-\lambda B_G(\theta_g \| \theta_d)) = p(\theta_g | \lambda\theta_d^*, \lambda) \exp(-\lambda G(\theta_d)). \qquad (23)$$

Here $\exp(-\lambda G(\theta_d))$ can be thought of as a prior on the discriminative parameters $p(\theta_d)$. In the above expression, $\exp(-\lambda B_G(\theta_g \| \theta_d)) = p(\theta_g | \theta_g) p(\theta_d)$ behaves as a joint coupling prior $P(\theta_d, \theta_g)$ as originally expected in the model (14). Note that hyperparameters of the prior $\alpha$ and $\beta$ are naturally derived from the geometric view of the conjugate prior. Here $\alpha = \lambda\theta_d^*$ and $\beta = \lambda$.

**Relation with Agarwal et al.:**

The prior we derived in the previous section turns out to be the exactly same as that proposed by Agarwal et al. [1], even though theirs was not formally justified. In that work, the authors break the coupled prior $p(\theta_g, \theta_d)$ into two parts: $p(\theta_d)$

and $p(\theta_g | \theta_d)$. They then derive an expression for the $p(\theta_g | \theta_d)$ based on the intuition that the mode of $p(\theta_g | \theta_d)$ should be $\theta_d$. Our analysis takes a different approach by coupling two models with the Bregman divergence rather than prior, and results in the expression and hyperparameters for the prior same as in [1].

## 7 Related Work and Conclusion

To our knowledge, there have been no previous attempts to understand Bayesian priors from a geometric perspective. One related piece of work [8] uses the Bayesian framework to find the best prior for a given distribution. It is noted that, in that work, the authors use the $\delta$-geometry for the data space and the $\alpha$-geometry for the prior space, and then show the different cases for different values $(\delta, \alpha)$. We emphasize that even though it is possible to use different geometry for the both spaces, it always makes more sense to use the same geometry. As mentioned in *remark* 1 in [8], useful cases are obtained only when we consider the same geometry.

We have shown that by considering the geometry induced by a likelihood function, the natural prior that results is exactly the conjugate prior. We have used this geometric understanding of conjugate prior to derive the coupling prior for the discriminative/generative hybrid model. Our derivation naturally gives us the expression and the hyperparameters of this coupling prior.

## References

[1] Agarwal, A., Daumé III, H.: Exponential family hybrid semi-supervised learning. In: In IJCAI. Pasadena, CA (2009)

[2] Amari, S.I., Nagaoka, H.: Methods of Information Geometry (Translations of Mathematical Monographs). American Mathematical Society (April 2001)

[3] Banerjee, A., Merugu, S., Dhillon, I.S., Ghosh, J.: Clustering with bregman divergences. Journal of Machine Learning Research 6 (October 2005)

[4] Bouchard, G.: Bias-variance tradeoff in hybrid generative-discriminative models. In: ICMLA '07. pp. 124–129. IEEE Computer Society, Washington, DC, USA (2007)

[5] Collins, M., Dasgupta, S., Schapire, R.E.: A generalization of principal component analysis to the exponential family. In: In NIPS 14. MIT Press (2001)

[6] Druck, G., Pal, C., McCallum, A., Zhu, X.: Semi-supervised classification with hybrid generative/discriminative methods. In: KDD '07. pp. 280–289. ACM, New York, NY, USA (2007)

[7] Lasserre, J.A., Bishop, C.M., Minka, T.P.: Principled hybrids of generative and discriminative models. In: CVPR '06. pp. 87–94. IEEE Computer Society, Washington, DC, USA (2006)

[8] Snoussi, H., Mohammad-Djafari, A.: Information geometry and prior selection (2002)

[9] Wainwright, M., Jordan, M.: Graphical models, exponential families, and variational inference. Tech. rep., University of California, Berkeley (2003)

[10] Wainwright, M.J., Jordan, M.I.: Graphical models, exponential families, and variational inference. Found. Trends Mach. Learn. 1(1-2), 1–305 (2008)