

BAYESIAN LEARNING FOR REGRESSION USING A DIRICHLET PRIOR PROBABILITY DISTRIBUTION

Paul Rademacher

U.S. Naval Research Laboratory
Radar Division
Washington, DC 20375, USA
paul.rademacher@nrl.navy.mil

Miloš Doroslovački

The George Washington University
Department of Electrical and Computer Engineering
Washington, DC 20052, USA
doroslov@gwu.edu

ABSTRACT

When taking a Bayesian approach to machine learning applications, the performance of the learned function strongly depends on how well the prior distribution selected by the designer matches the true data-generating distribution. Dirichlet priors have a number of desirable properties - they lead to closed-form posterior distributions given independent training data, have full support over the space of data probability distributions, and can be fully objective or subjective depending on their localization parameter. This paper assumes a Dirichlet prior and produces predictive distributions to characterize unobservable random quantities given observed data. The results are then applied to the most common loss function for regression, the squared-error loss. The optimal Bayes estimator and the corresponding minimum risk are presented and interpreted for different values of the localization parameter; specific attention is given to the extremal values.

Index Terms— Bayesian learning, machine learning, regression, estimation, Dirichlet distribution, predictive distribution

1. INTRODUCTION

The success or failure of Bayesian learning methods hinge on how well the prior knowledge imparted by the designer matches reality. The chosen prior distribution over the set of data-generating probability distributions reflects the users confidence that different distributions are responsible for generating the observed/unobserved random elements. If a highly localized prior is chosen that strongly weights the actual data probability distribution, low risk learning functions are possible even with limited training data; however, if the localized prior is poorly selected, a good solution is unlikely to be achieved. Conversely, a less localized prior that treats the different distributions without preference adapts more responsively during training; if data is limited, however, the learning function may not deliver the required performance.

This work assumes that the joint observed and unobserved data elements are drawn from finite sets. The probability

mass function (PMF) generating the data is characterized by a Dirichlet prior. The class of Dirichlet probability density functions (PDF) has the desirable properties of full support over the set of possible data-generating distributions and a closed-form posterior distribution for independently and identically distributed data [1]. Furthermore, control of the Dirichlet parameters can enable both minimally and maximally localized priors; the objective uniform distribution and the subjective delta function distribution are special cases.

After introducing the problem and discussing the relevant probability distributions, the Bayesian framework will be applied to the most common loss function for regression: the squared error loss function. Algebraic forms of optimal estimators and their corresponding minimum risk will be presented. Specific attention will be given to performance for the asymptotic cases of Dirichlet prior localization.

2. OBJECTIVE

Consider an observable random element $x \in \mathcal{X}$ and an unobservable random element $y \in \mathcal{Y}$ which are jointly distributed according to an unknown probability distribution $\theta \in \Theta = \left\{ \theta \in \mathbb{R}_{\geq 0}^{\mathcal{Y} \times \mathcal{X}} : \sum_{y \in \mathcal{Y}} \sum_{x \in \mathcal{X}} \theta(y, x) = 1 \right\}$, such that $P_{y,x|\theta}(y, x|\theta) = \theta(y, x)$.

Also observed is a random sequence of N samples drawn from θ , denoted $D = (Y, X) \in \mathcal{D} = \{\mathcal{Y} \times \mathcal{X}\}^N$. The N data pairs are identically distributed as $P_{D_n|\theta}(y, x|\theta) = \theta(y, x)$ and are conditionally independent from one another and from the novel pair (y, x) .

The objective is to design a decision function $f : \mathcal{D} \mapsto \mathcal{H}^{\mathcal{X}}$ which outputs a mapping from the space of the observed random element x to a decision space \mathcal{H} . The metric guiding the design is a loss function $\mathcal{L} : \mathcal{H} \times \mathcal{Y} \mapsto \mathbb{R}_{\geq 0}$ which penalizes the decision $h \in \mathcal{H}$ based on the value of y . The conditional expected loss, or conditional “risk”, is defined as

$$\mathcal{R}_{\Theta}(f; \theta) = E_{D|\theta} \left[E_{y,x|\theta} \left[\mathcal{L}(f(x; D), y) \right] \right]. \quad (1)$$

As the model θ is not observed, \mathcal{R}_{Θ} is not yet a feasible

objective function for optimization. If the designer selects a PDF p_θ , the Bayes risk is formulated as

$$\begin{aligned} \mathcal{R}(f) &= \mathbb{E}_\theta [\mathcal{R}_\Theta(f; \theta)] \\ &= \mathbb{E}_{y, x, D} [\mathcal{L}(f(x; D), y)] \end{aligned} \quad (2)$$

and y , x , and D are treated as jointly distributed random elements. The optimal learning function is expressed as

$$f^*(x; D) = \arg \min_{h \in \mathcal{H}} \mathbb{E}_{y | x, D} [\mathcal{L}(h, y)] \quad (3)$$

and the corresponding minimum Bayes risk is

$$\mathcal{R}(f^*) = \mathbb{E}_{x, D} \left[\min_{h \in \mathcal{H}} \mathbb{E}_{y | x, D} [\mathcal{L}(h, y)] \right]. \quad (4)$$

3. PROBABILITY DISTRIBUTIONS

In this section, the joint PMF $P_{y, x, D}$ is determined using the Dirichlet prior p_θ . Other distributions of interest will be provided, including the training data PMF P_D and the predictive distribution $P_{y | x, D}$.

3.1. Model PDF, p_θ

The Dirichlet PDF for the model random process $\theta \in \Theta$ is [2]

$$p_\theta(\theta) = \beta(\alpha)^{-1} \prod_{y \in \mathcal{Y}} \prod_{x \in \mathcal{X}} \theta(y, x)^{\alpha(y, x) - 1}, \quad (5)$$

where the user-selected PDF parameters $\alpha : \mathcal{Y} \times \mathcal{X} \mapsto \mathbb{R}^+$ are introduced and β is the multivariate beta function.

The parameters α controls around which models θ the PDF concentrates and how strongly. For convenience, introduce the localization parameter $\alpha_0 \equiv \sum_{y \in \mathcal{Y}} \sum_{x \in \mathcal{X}} \alpha(y, x)$. Note that $P_{y, x}(y, x) = \mu_\theta(y, x) = \alpha(y, x) / \alpha_0$, where μ denotes the mean function of a random variable.

Of specific interest is how p_θ changes as the localization parameter approaches its limiting values. For $\alpha_0 \rightarrow \infty$, the PDF concentrates at its mean, resulting in

$$p_\theta(\theta) \rightarrow \delta\left(\theta - \frac{\alpha}{\alpha_0}\right), \quad (6)$$

where $\delta(\cdot)$ represents the Dirac delta function over Θ . Conversely, for $\alpha_0 \rightarrow 0$, the PDF tends toward

$$p_\theta(\theta) \rightarrow \sum_{y \in \mathcal{Y}} \sum_{x \in \mathcal{X}} \frac{\alpha(y, x)}{\alpha_0} \delta(\theta(\cdot, \cdot) - \delta[\cdot, y] \delta[\cdot, x]), \quad (7)$$

where $\delta[\cdot, \cdot]$ is the Kronecker delta function. Note that the PDF has support only for the $|\mathcal{Y}| |\mathcal{X}|$ models with an ℓ_0 norm satisfying $\|\theta\|_0 = 1$.

These trends are demonstrated in Figure 1. The cardinalities $|\mathcal{Y}| = 3$ and $|\mathcal{X}| = 1$ are chosen to enable visualization, despite the implication that x is deterministic. Note that for the top sub-plot, $\alpha(y, x) < 1$ and the PDF tends to infinity at the domain boundaries; this cannot be captured by the plot color scale.

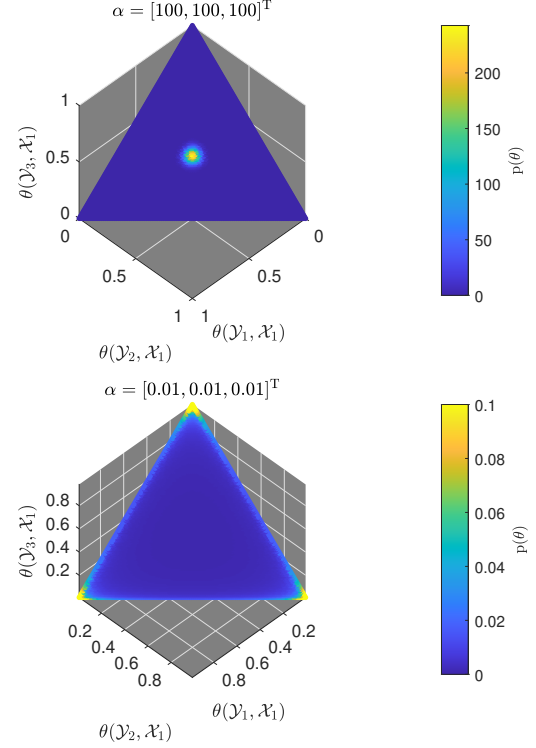


Fig. 1. Model prior PDF p_θ for different localizations α_0

3.2. Training Set PMF, P_D

The distribution of the training data D conditioned on the model can be formulated as

$$P_{D | \theta}(D | \theta) = \prod_{y \in \mathcal{Y}} \prod_{x \in \mathcal{X}} \theta(y, x)^{\bar{N}(y, x; D)}, \quad (8)$$

where the dependency on the training data D is expressed through a transform function $\bar{N} : \mathcal{D} \mapsto \bar{\mathcal{N}}$ defined as $\bar{N}(y, x; D) = \sum_{n=1}^N \delta[y, Y_n] \delta[x, X_n]$. This transform counts the occurrences of each possible pair (y, x) in the training data; its range is a function space $\bar{\mathcal{N}} \subset \mathbb{Z}_{\geq 0}^{\mathcal{Y} \times \mathcal{X}}$.

As the conditional distribution $P_{D | \theta}$ is of exponential form, it can be readily shown that the marginal distribution of the training data is [3]

$$P_D(D) = \beta(\alpha)^{-1} \beta(\alpha + \bar{N}(D)). \quad (9)$$

Note that both $P_{D | \theta}$ and P_D depend on the training data D only through the transform \bar{N} ; consequently, $\bar{N}(D)$ is a sufficient statistic [4] for the model θ . As such, it is useful to define a new random process $\bar{n} \equiv \bar{N}(D) \in \bar{\mathcal{N}}$.

The cardinality of the random process' domain is $|\bar{\mathcal{N}}| = \mathcal{M}(\{N, |\mathcal{Y}| |\mathcal{X}| - 1\})$, where \mathcal{M} is the multinomial coefficient; this can be shown using the stars-and-bars method [5]. The cardinality of original set is $|\mathcal{D}| = (|\mathcal{Y}| |\mathcal{X}|)^N$; thus $|\bar{\mathcal{N}}| \leq |\mathcal{D}|$ and the sufficient statistic compactly represents the valuable information in the training data.

It is easily shown that the conditional PMF $P_{\bar{n}|\theta}$ is multinomial. As the Dirichlet distribution characterizes the parameters of this multinomial distribution, the marginal PMF of \bar{n} is a Dirichlet-Multinomial distribution [6] parameterized by α ,

$$P_{\bar{n}}(\bar{n}) = \mathcal{M}(\bar{n})\beta(\alpha)^{-1}\beta(\alpha + \bar{n}). \quad (10)$$

The first and second joint moments of \bar{n} are

$$\mu_{\bar{n}}(y, x) = N \frac{\alpha(y, x)}{\alpha_0} = N\mu_{\theta}(y, x) \quad (11)$$

and

$$\begin{aligned} E[\bar{n}(y, x)\bar{n}(y', x')] & \quad (12) \\ &= \frac{N}{\alpha_0 + 1} \left((\alpha_0 + N)\mu_{\theta}(y, x)\delta[y, y']\delta[x, x'] \right. \\ & \quad \left. + \alpha_0(N - 1)\mu_{\theta}(y, x)\mu_{\theta}(y', x') \right). \end{aligned}$$

Again, the distributions for minimal and maximal α_0 are relevant. For $\alpha_0 \rightarrow \infty$, the model PDF p_{θ} concentrates at its mean and thus \bar{n} is characterized by a multinomial distribution,

$$P_{\bar{n}}(\bar{n}) = \mathcal{M}(\bar{n}) \prod_{y \in \mathcal{Y}} \prod_{x \in \mathcal{X}} \left(\frac{\alpha(y, x)}{\alpha_0} \right)^{\bar{n}(y, x)}. \quad (13)$$

Conversely, for $\alpha_0 \rightarrow 0$, the PMF tends toward

$$P_{\bar{n}}(\bar{n}) = \sum_{y \in \mathcal{Y}} \sum_{x \in \mathcal{X}} \frac{\alpha(y, x)}{\alpha_0} \delta[\bar{n}(\cdot, \cdot), N\delta[\cdot, y]\delta[\cdot, x]] \quad (14)$$

and the training data are identical.

3.2.1. Aggregation Properties

Also of importance is the “marginalized” random process n' , defined as $n'(x) \equiv \sum_{y \in \mathcal{Y}} \bar{n}(y, x)$ over the set \mathcal{X} . By the aggregation property of Dirichlet-Multinomial functions [6], the new function is distributed as $n' \sim \text{DM}(N, \alpha')$, where $\alpha'(x) = \sum_{y \in \mathcal{Y}} \alpha(y, x)$.

Also of interest is the distribution of \bar{n} conditioned on its aggregation n' . Using the Dirichlet-Multinomial properties presented in Appendix A,

$$\begin{aligned} P_{\bar{n}|n'}(\bar{n}|n') & \quad (15) \\ &= \prod_{x \in \mathcal{X}} \left[\mathcal{M}(\bar{n}(\cdot, x))\beta(\alpha(\cdot, x))^{-1}\beta(\alpha(\cdot, x) + \bar{n}(\cdot, x)) \right]. \end{aligned}$$

Observe that conditioning on the aggregation renders the function segments $\bar{n}(\cdot, x)$ independent of one another and that they are also Dirichlet-Multinomial.

3.3. Predictive PMF, $P_{y|x,D}$

As shown in Equation (3), the decision selected by the optimally designed function depends on the predictive distribution of the unobserved y conditioned on all observable random elements. Note that $P_{y|x,D} = E_{\theta|x,D} [P_{y|x,\theta}]$ - the Bayesian predictive PMF is the expected value of the true predictive PMF $P_{y|x,\theta}$ with respect to the model posterior distribution.

Since $P_{D|\theta}$ is of exponential form, the Dirichlet prior p_{θ} is its conjugate prior [7]; thus, the model posterior PDF given the training data is

$$\begin{aligned} P_{\theta|D}(\theta|D) &= \frac{P_{D|\theta}(D|\theta) p_{\theta}(\theta)}{P_D(D)} \quad (16) \\ &= \beta(\alpha + \bar{N}(D))^{-1} \prod_{y \in \mathcal{Y}} \prod_{x \in \mathcal{X}} \theta(y, x)^{\alpha(y, x) + \bar{N}(y, x; D) - 1}, \end{aligned}$$

a Dirichlet distribution with parameter function $\alpha + \bar{N}(D)$.

Also, the localization parameter increases proportionately with the volume of training data; consequently, as $N \rightarrow \infty$, the posterior converges to $p_{\theta|D}(\theta|D) \rightarrow \delta(\theta - \bar{N}(D)/N)$ and the model is positively identified. This is a consequence of the full support of the Dirichlet distribution; general posterior distributions do not tend to a delta function at the empirical PMF. Conversely, as $\alpha_0 \rightarrow \infty$, the prior model certainty is stronger and the posterior tends toward $p_{\theta|D}(\theta|D) \rightarrow \delta(\theta - \alpha/\alpha_0)$, independent of the training data.

The joint PMF of y and x conditioned on the training data is expressed as [8]

$$\begin{aligned} P_{y,x|D}(y, x|D) &= \mu_{\theta|D}(y, x) \quad (17) \\ &= \frac{\alpha(y, x) + \bar{N}(y, x; D)}{\alpha_0 + N} \end{aligned}$$

and the predictive distribution is generated via Bayes rule as

$$\begin{aligned} P_{y|x,D}(y|x, D) &= \frac{\alpha(y, x) + \bar{N}(y, x; D)}{\alpha'(x) + N'(x; D)} \quad (18) \\ &= \left(\frac{\alpha'(x)}{\alpha'(x) + N'(x; D)} \right) P_{y|x}(y|x) \\ & \quad + \left(\frac{N'(x; D)}{\alpha'(x) + N'(x; D)} \right) \frac{\bar{N}(y, x; D)}{N'(x; D)}, \end{aligned}$$

where $N'(x; D) = \sum_{n=1}^N \delta[x, X_n]$. The last representation views the distribution as a convex combination of two conditional distributions. The first distribution $P_{y|x}(y|x) = \alpha(y, x)/\alpha'(x)$ is independent of the training data and based on the prior knowledge implied via the model PDF parameter; the second distribution is the conditional empirical PMF and depends solely on D . For both, only those values α and D corresponding to the observed value of x influence the distribution.

The weighting factors are dependent on these values as well. For $N'(x; D) = 0$ or as $\alpha_0 \rightarrow \infty$, the PMF tends toward the conditional distribution $P_{y|x}$, which only depends

on the model parameter α . As the number of training examples increases or as $\alpha_0 \rightarrow 0$, $P_{y|x,D}$ tends towards the empirical conditional distribution.

4. REGRESSION AND THE SQUARED-ERROR LOSS

The squared-error (SE) loss function is arguably the most commonly used loss function for regression, or in fact for any estimation problem. This can be attributed to its quadratic form, which enables closed-form expression of the minimizing estimation function f^* .

It is assumed that the unobserved random element y is a scalar random variable; that is, $\mathcal{Y} \subset \mathbb{R}$. Additionally, the learning function's estimate is allowed to assume real numbers; thus, $\mathcal{H} = \mathbb{R} \supset \mathcal{Y}$.

The loss function is defined as

$$\mathcal{L}(h, y) = (h - y)^2. \quad (19)$$

Substituting into (2), the Bayes risk for a general estimator is

$$\mathcal{R}(f) = E_{x,D} \left[E_{y|x,D} \left[(f(x; D) - y)^2 \right] \right]. \quad (20)$$

4.1. Optimal Estimate: the Posterior Mean

To find the optimal estimator, the squared-error loss is substituted into (3); note that the objective function is quadratic over the argument h . It is easily shown that the function over h is positive-definite; as such, the minimizing decision h is the sole stationary point. Setting the first derivative of the function to zero, the optimal estimate is the expected value of y given the training data and the observed value x , such that

$$\begin{aligned} f^*(x; D) &= \arg \min_{h \in \mathbb{R}} E_{y|x,D} [(h - y)^2] = \mu_{y|x,D} \quad (21) \\ &= \left(\frac{\alpha'(x)}{\alpha'(x) + N'(x; D)} \right) \mu_{y|x} \\ &\quad + \left(\frac{N'(x; D)}{\alpha'(x) + N'(x; D)} \right) \frac{\sum_{n=1}^N \delta[x, X_n] Y_n}{N'(x; D)}. \end{aligned}$$

The optimal estimate is interpreted as a convex combination of two separate estimates - the expected value of y conditioned on the observed x and the average of the training values Y_n which have a value X_n matching the observed value of x . The weighting factors are the same as those of $P_{y|x,D}$; thus, stronger prior information (larger $\alpha'(x)$) provides more weight to the estimate $\mu_{y|x}$ and more voluminous training data puts emphasis on the empirical conditional mean.

Another interesting form for the optimal estimator is $f^*(x; D) = E_{\theta|x,D} [\mu_{y|x,\theta}]$. If the model θ were known, then the clairvoyant estimate $\mu_{y|x,\theta}$ would be optimal; instead, all such estimates are weighted and combined via the expectation over the model posterior PDF $p_{\theta|x,D}$.

4.2. Minimum Risk: the Expected Posterior Variance

Substituting the optimal estimator (21) into Equation (20), the minimum Bayes risk is the expected conditional variance

$$\begin{aligned} \mathcal{R}(f^*) &= E_{x,D} [\Sigma_{y|x,D}] \quad (22) \\ &= E_{x,\theta} [\Sigma_{y|x,\theta}] + E_{x,D} [C_{\theta|x,D} [\mu_{y|x,\theta}]], \end{aligned}$$

where Σ is the variance of a random variable and the operator C is the variance of a function of a random variable. The second formula is of interest. The first term is the expected squared-error of the clairvoyant estimate $\mu_{y|x,\theta}$; the second term is its expected conditional variance with respect to the model posterior PDF $p_{\theta|x,D}$.

Using the sufficient statistic $\bar{n} \equiv \bar{N}(D)$, the minimum risk can also be represented as $E_{x,\bar{n}} [\Sigma_{y|x,\bar{n}}]$. Decompose the conditional variance as

$$\Sigma_{y|x,\bar{n}} = E_{y|x,\bar{n}} [y^2] - \mu_{y|x,\bar{n}}^2 \quad (23)$$

and assess the expected values of these terms separately. The first term is

$$E_{x,\bar{n}} [E_{y|x,\bar{n}} [y^2]] = E_x [E_{y|x} [y^2]], \quad (24)$$

where the expectation over x is performed using $P_x(x) = \alpha'(x)/\alpha_0$. As proven in Appendix B, the second term is

$$\begin{aligned} E_{x,\bar{n}} [\mu_{y|x,\bar{n}}^2] \quad (25) \\ = E_x \left[\frac{N E_{y|x} [y^2] + (\alpha_0 \alpha'(x) + N \alpha'(x) + \alpha_0) \mu_{y|x}^2}{(\alpha_0 + N)(\alpha'(x) + 1)} \right]. \end{aligned}$$

Combining, the minimum Bayes risk is

$$\begin{aligned} \mathcal{R}(f^*) &= E_{x,\bar{n}} [E_{y|x,\bar{n}} [y^2] - \mu_{y|x,\bar{n}}^2] \quad (26) \\ &= E_x \left[\frac{P_x(x) + (\alpha_0 + N)^{-1}}{P_x(x) + \alpha_0^{-1}} \Sigma_{y|x} \right]. \end{aligned}$$

The risk is a convex combination of scaled conditional variances for the different PMFs $P_{y|x}$. The convex coefficients are values from the prior marginal distribution P_x .

The scaling factor for each term $\Sigma_{y|x}$ depends on the marginal PMF value $P_x(x)$, as well as on the prior localization α_0 and the number of training samples N . Observe that with no training data ($N = 0$), the scaling factor becomes unity and the risk is $\mathcal{R}(f^*) = E_x [\Sigma_{y|x}]$. Conversely, as $N \rightarrow \infty$, the Bayes risk is $\mathcal{R}(f^*) = E_x \left[\frac{P_x(x)}{P_x(x) + \alpha_0^{-1}} \Sigma_{y|x} \right] = E_{x,\theta} [\Sigma_{y|x,\theta}]$. Also, as the model localization parameter $\alpha_0 \rightarrow 0$, the risk tends to zero (for $N > 0$); as $\alpha_0 \rightarrow \infty$, the risk tends toward $E_x [\Sigma_{y|x}]$.

4.2.1. Examples

To illustrate these trends, explicitly define the sets $\mathcal{Y} = \{i/M_y : i = 0, \dots, M_y - 1\}$ and $\mathcal{X} = \{i/M_x : i =$

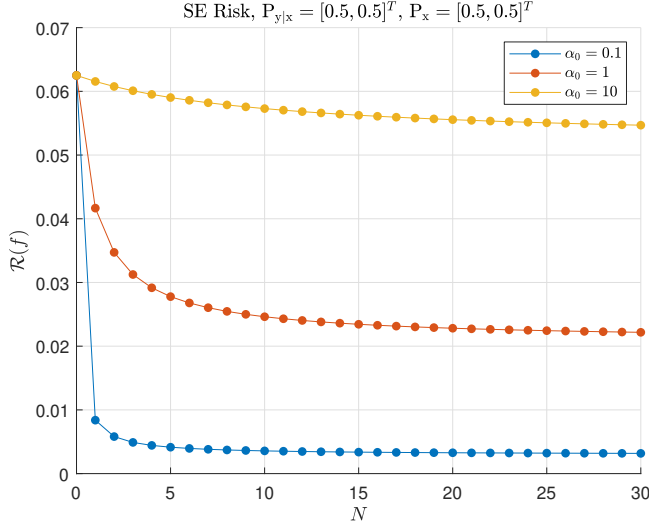


Fig. 2. Minimum SE Risk for different training set sizes N

$0, \dots, M_x - 1\}$. Let the conditional variances $\Sigma_{y|x}$ be the same for each value $x = x_i$; in this case, the squared-error becomes the conditional variance scaled by a factor dependent on the marginal distribution P_x , such that $\mathcal{R}(f^*) = \Sigma_{y|x} E_x \left[\frac{P_x(x) + (\alpha_0 + N)^{-1}}{P_x(x) + \alpha_0^{-1}} \right]$. Figures 2 and 3 display how the risk changes with N and α_0 when $P_{y|x}$ and P_x are fixed.

It may not seem intuitive for the risk to decrease when α_0 is smaller – the variance of the model θ increases and the prior knowledge is less definitive. This is a result of the Dirichlet PDF weight shifting towards the $|\mathcal{Y}| |\mathcal{X}|$ different sparse models which have ℓ_0 norms satisfying $\|\theta\|_0 = 1$. Although these PMFs are maximally separated, they all have zero variance. The optimal learner (21) will simply use the empirical distribution supplied via the training data - this allows exact identification of θ with a single training pair.

It is also informative to visualize how the minimum squared-error changes with the marginal distribution P_x for fixed volume of training data N and prior localization α_0 . Figure 4 demonstrates how the risk changes with this marginal PMF. Observe that the risk is maximal at the distributions satisfying $\|P_x\|_0 = 1$; the scaling factor for the conditional variance $\Sigma_{y|x}$ becomes $\frac{1 + (\alpha_0 + N)^{-1}}{1 + \alpha_0^{-1}}$. Conversely, for $P_x = 1/|\mathcal{X}|$ the scaling factor becomes $\frac{|\mathcal{X}|^{-1} + (\alpha_0 + N)^{-1}}{|\mathcal{X}|^{-1} + \alpha_0^{-1}}$ and the risk is minimal.

5. CONCLUSIONS

This paper has assumed a Dirichlet prior for Bayesian learning, established relevant distributions, and applied the framework to squared-error regression. Closed-forms have been provided for the optimal estimation function and the mini-

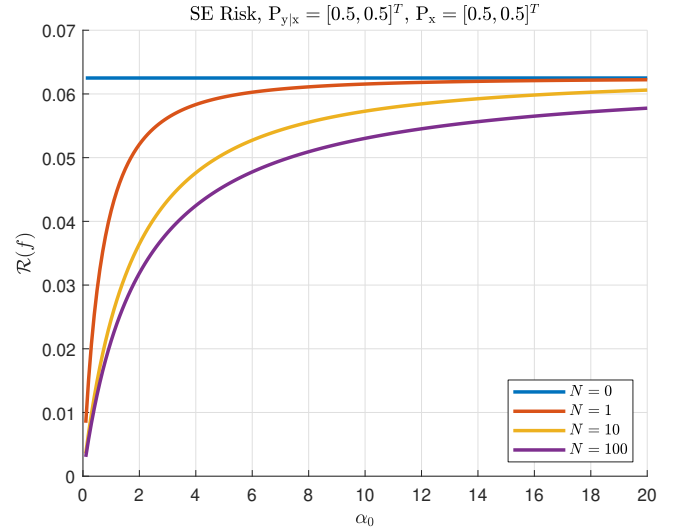


Fig. 3. Minimum SE Risk for different prior localization α_0

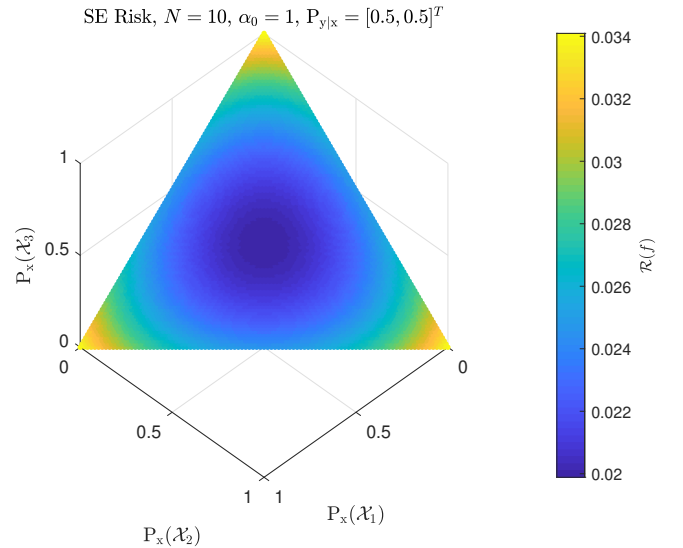


Fig. 4. Minimum SE Risk for different PMFs P_x

mum Bayes risk. Analysis and graphical examples highlight interesting trends for the estimation squared-error as a function of training data volume and the Dirichlet prior distribution parameters. Notably, Dirichlet priors with smaller localization parameters (and thus larger covariance) lead to optimal estimators with lower Bayes risk.

Future work will explore how the conditional risk (1) changes with different estimator parameterizations α for different true models θ . It will be shown that there is an optimal localization parameter dependent on the Dirichlet prior mean μ_θ and the actual model; whether the optimal localization parameter is large or small will depend on the quality of the match between μ_θ and θ . Additionally, a “minimax” function will be found that minimizes the worst-case conditional risk $\max_{\theta \in \Theta} \mathcal{R}_\Theta(f; \theta)$.

A. DIRICHLET-MULTINOMIAL RANDOM PROCESS CONDITIONED ON ITS AGGREGATION

A defining characteristic of a Dirichlet-Multinomial random process is that its aggregations are also Dirichlet-Multinomial [6]. Consider a DM random process $\bar{n} \sim \text{DM}(N, \alpha)$ over the set \mathcal{Y} . Define an arbitrary partition of \mathcal{Y} : $\{\dots, \mathcal{S}(z), \dots\}$, $z \in \mathcal{Z}$; the transformed random process $n'(z) \equiv \sum_{y \in \mathcal{S}(z)} \bar{n}(y)$ is necessarily Dirichlet-Multinomial with parameterizing function $\alpha'(z) = \sum_{y \in \mathcal{S}(z)} \alpha(y)$.

Conditioned on the aggregation n' , the segments $\{\bar{n}(y) : y \in \mathcal{S}(z)\}$ of the original random process become independent Dirichlet-Multinomial random processes, such that

$$\begin{aligned} P_{\bar{n}|n'}(\bar{n}|n') &= \frac{P_{\bar{n}}(\bar{n})}{P_{n'}(n')} P_{n'|\bar{n}}(n'|\bar{n}) \\ &= \frac{\mathcal{M}(\bar{n})\beta(\alpha)^{-1}\beta(\alpha + \bar{n})}{\mathcal{M}(n')\beta(\alpha')^{-1}\beta(\alpha' + n')} \\ &= \left(\prod_{z \in \mathcal{Z}} \frac{\Gamma(\alpha'(z) + n'(z))}{n'(z)!\Gamma(\alpha'(z))} \right)^{-1} \left(\prod_{y \in \mathcal{Y}} \frac{\Gamma(\alpha(y) + \bar{n}(y))}{\bar{n}(y)!\Gamma(\alpha(y))} \right) \\ &= \prod_{z \in \mathcal{Z}} \left[\frac{n'(z)!\Gamma(\alpha'(z))}{\Gamma(\alpha'(z) + n'(z))} \prod_{y \in \mathcal{S}(z)} \frac{\Gamma(\alpha(y) + \bar{n}(y))}{\bar{n}(y)!\Gamma(\alpha(y))} \right], \end{aligned} \quad (27)$$

over domain $\bar{n} \in \left\{ \mathbb{Z}_{\geq 0}^{\mathcal{Y}} : \sum_{y \in \mathcal{S}(z)} \bar{n}(y) = n'(z), \quad \forall z \in \mathcal{Z} \right\}$.

B. FORMULATION OF $E_{x, \bar{n}} \left[\mu_{y|x, \bar{n}}^2 \right]$

The expected value of the conditional squared mean is compactly represented as

$$\begin{aligned} E_{x, \bar{n}} \left[\mu_{y|x, \bar{n}}^2 \right] \\ = \sum_{x, y, y'} yy' E_{\bar{n}} \left[P_{y, x|\bar{n}}(y, x|\bar{n}) P_{y|x, \bar{n}}(y'|x, \bar{n}) \right] \end{aligned} \quad (28)$$

$$\begin{aligned} &= \sum_{x, y, y'} yy' E_{n'} \left[\frac{1}{(\alpha_0 + N)(\alpha'(x) + n'(x))} \right. \\ &\quad \left. E_{\bar{n}|n'} \left[(\alpha(y, x) + \bar{n}(y, x))(\alpha(y', x) + \bar{n}(y', x)) \right] \right] \\ &= \sum_{x \in \mathcal{X}} \frac{1}{(\alpha_0 + N)(\alpha'(x) + 1)} E_{n'} \left[n'(x) E_{y|x}[y^2](x) \right. \\ &\quad \left. + \alpha'(x)(\alpha'(x) + n'(x) + 1)\mu_{y|x}^2(x) \right] \\ &= E_x \left[\frac{N E_{y|x}[y^2] + (\alpha_0 \alpha'(x) + N \alpha'(x) + \alpha_0)\mu_{y|x}^2}{(\alpha_0 + N)(\alpha'(x) + 1)} \right]. \end{aligned}$$

The formulation uses the statistical characterization of the aggregation, $n' \sim \text{DM}(N, \alpha')$; also used is the property that the Dirichlet-Multinomial random process \bar{n} conditioned on its aggregation n' yields independent conditional DM functions $\bar{n}(\cdot, x) | n'(x) \sim \text{DM}(n'(x), \alpha(\cdot, x))$.

C. REFERENCES

- [1] T. S. Ferguson, “A bayesian analysis of some nonparametric problems,” *The Annals of Statistics*, vol. 1, no. 2, pp. 209–230, 1973.
- [2] C. M. Bishop, *Pattern Recognition and Machine Learning*, ser. Information Science and Statistics. Springer, 2006.
- [3] T. P. Minka, “Bayesian inference, entropy, and the multinomial distribution,” Microsoft Research, Tech. Rep., 2003.
- [4] J. M. Bernardo and A. F. M. Smith, *Bayesian Theory*, ser. Wiley Series in Probability and Statistics. John Wiley & Sons, 2000.
- [5] W. Feller, *An Introduction to Probability Theory and Its Applications*, 2nd ed., ser. Probability and Mathematical Statistics. New York, New York: John Wiley & Sons, 1971, vol. 2.
- [6] N. L. Johnson, S. Kotz, and N. Balakrishnan, *Discrete Multivariate Distributions*, ser. Probability and Statistics. John Wiley & Sons, 1997.
- [7] S. Theodoridis, *Machine Learning: A Bayesian and Optimization Perspective*. Elsevier, 2015.
- [8] K. P. Murphy, “Binomial and multinomial distributions,” University of British Columbia, Tech. Rep., 2006.