

# Bayesian Learning using a Dirichlet Prior for Regression and Classification

Paul Rademacher

January 17, 2020

# Contents

<b>1</b>	<b>Introduction</b>	<b>8</b>
1.1	Background . . . . .	8
1.2	Notation . . . . .	9
<b>2</b>	<b>Problem Statement</b>	<b>14</b>
2.1	Data Model and Objective . . . . .	14
2.1.1	Clairvoyant Decision . . . . .	15
2.1.2	Bayes Decision . . . . .	16
2.1.2.1	Irreducible Risk . . . . .	17
2.2	Sufficient Statistic: the Empirical PMF . . . . .	17
2.3	Marginal/Conditional Distributions for Observed/Unobserved Elements . . .	19
2.3.1	Marginal and Conditional Distributions of $\theta$ . . . . .	20
2.3.2	Marginal and Conditional Distributions of Training Data . . . . .	20
2.3.3	Model Posteriors . . . . .	21
2.4	Applications to Common Loss Functions . . . . .	22
2.4.1	Regression: the Squared-Error Loss . . . . .	23
2.4.1.1	Clairvoyant Estimation . . . . .	23
2.4.1.2	Bayesian Estimation . . . . .	24
2.4.2	Classification: the 0-1 Loss . . . . .	26
2.4.2.1	Clairvoyant Hypothesis . . . . .	26
2.4.2.2	Bayesian Classification . . . . .	27

<b>3 Finite Dirichlet Model</b>	<b>29</b>
3.1 Probability Distributions . . . . .	29
3.1.1 Model PDF, $p_\theta$ . . . . .	29
3.1.1.1 Marginal and Conditional Distributions . . . . .	31
3.1.2 Training Set PMF, $P_D$ . . . . .	34
3.1.2.1 Marginal and Conditional Distributions . . . . .	39
3.1.3 Predictive PMF, $P_{y x,D}$ . . . . .	40
3.1.3.1 Representation using the complete model posterior . . . . .	44
3.2 Model Estimation Perspective . . . . .	45
3.3 Applications to Common Loss Functions . . . . .	49
3.3.1 Regression: the Squared-Error Loss . . . . .	52
3.3.1.1 Optimal Estimate: the Posterior Mean . . . . .	52
3.3.1.2 Minimum Risk: the Expected Posterior Variance . . . . .	53
3.3.1.3 Conditional Squared-Error for a Dirichlet-based Estimator .	60
3.3.2 Classification: the 0-1 Loss . . . . .	66
3.3.2.1 Optimal Hypothesis: Conditional Maximum <i>a posteriori</i> .	66
3.3.2.2 Minimum Risk: Probability of Error . . . . .	68
3.3.2.3 Conditional Probability of Error for a Dirichlet-based Classifier	76
<b>4 Extention to Infinite-Dimensional Spaces - Countably Infinite</b>	<b>80</b>
4.1 Intro . . . . .	80
4.2 Basic Model . . . . .	80
4.2.1 Probability Distributions . . . . .	80
4.2.1.1 Model PDF, $p(\theta)$ . . . . .	80
4.2.1.2 Training Data PMF, $P(D)$ . . . . .	81
4.2.1.3 Output conditional PMF, $P(y D)$ . . . . .	82
4.3 Application to Common Loss Functions . . . . .	82
4.4 General Model . . . . .	83
4.5 Applications: General Model . . . . .	83

<b>5 Extention to Infinite-Dimensional Spaces - Uncountably Infinite</b>	<b>84</b>
5.1 Intro . . . . .	84
5.2 Basic Model . . . . .	84
5.2.1 Probability Distributions . . . . .	84
5.2.1.1 Model $\theta$ Characterization . . . . .	84
5.2.1.2 Output conditional PDF, $p_{y D}$ . . . . .	85
5.2.1.3 Training Data PDF, $p(D)$ . . . . .	85
5.3 Application to Common Loss Functions . . . . .	87
5.3.1 Regression: the Squared-Error Loss . . . . .	87
5.3.1.1 Optimal Learner . . . . .	88
5.3.1.2 Minimum Risk . . . . .	88
5.4 General Model . . . . .	89
5.4.1 Model Extension . . . . .	90
5.4.2 General Probability Distributions . . . . .	90
5.4.2.1 Model $\theta$ Characterization . . . . .	90
5.4.2.2 Output conditional PDF, $p_{y x,D}$ . . . . .	90
5.4.2.3 Training Data PDF, $p_D$ . . . . .	91
5.5 Applications: General Model . . . . .	94
5.5.1 Regression: the Squared-Error Loss . . . . .	94
5.5.1.1 Optimal Learner . . . . .	95
5.5.1.2 Minimum Risk . . . . .	95
<b>A</b>	<b>99</b>
A.1 Dirichlet random process conditioned on its aggregation . . . . .	99
A.2 Multinomial Distribution Properties . . . . .	100
A.2.1 Aggregation . . . . .	100
A.2.2 Conditioned on its Aggregation . . . . .	101
A.3 Dirichlet-Multinomial random process conditioned on its aggregation . . . . .	101
A.4 First and Second moments of a Dirichlet Process . . . . .	102
A.5 Proof: Continuous Model Posterior Distribution is Dirichlet Process . . . . .	103

A.6 The Dirichlet-Multinomial Process . . . . .	105
A.6.1 Definition . . . . .	105
A.6.2 Proof that $\sum_{n=1}^N \delta(y - D_n)$ is a DMP . . . . .	105
A.6.3 Mean and Correlation Functions . . . . .	106
A.6.4 Continuous aggregation . . . . .	107
<b>B</b>	<b>109</b>
B.1 Maximum <i>a Posteriori</i> estimate of $\theta$ given $D$ . . . . .	109
<b>Bibliography</b>	<b>110</b>

# Todo list

█ use todonotes package instead of my initials? . . . . .	6
█ equation numbers to final line! . . . . .	6
█ brackets for expectation ops? . . . . .	6
█ line break symbol format, before/after? . . . . .	6
█ learners (full) or decision functions (range)? . . . . .	6
█ likelihood function terminology . . . . .	6
█ Dirichlet localization or concentration? . . . . .	6
█ fix arguments for theta, nbar dists. . . . .	6
█ suppress arguments where sensible? eg Pd = theta . . . . .	6
█ ALL figure notation: theta font + Ycal indexing. use R,f opt? . . . . .	6
█ is $\Theta$ redundant given $\mathcal{P}$ ? . . . . .	6
█ DIM and PR operator from AMS? . . . . .	6
█ multinomial random process? . . . . .	6
█ introduce tilde alpha to match tilde theta? . . . . .	6
█ use x,y subs instead of prime/tilde? . . . . .	6
█ change from nbar to empirical RV??? . . . . .	6
█ investigate N lim for nbar given theta, bayes risk vs model support . . . . .	6
█ SEMI-SUPERVISED - generalize to joint decisions!!! training/test! . . . . .	6
█ priors = sparse conditionals; w/ sufficient statistics . . . . .	6
█ NFLT investigation? try matlab examples . . . . .	6
█ use clairvoyant/irreducible terms and symbols . . . . .	6
█ empirical risk terms/discussion? . . . . .	6
█ check notation throughout for expectation/variance operators, proper arguments .	7

generalize y,x,h from scalars to functions!!! . . . . .	7
jeffrey prior, fisher info? . . . . .	7
aleph reference? . . . . .	10
explicit PMF/PDF formula with P of events? . . . . .	10
Ever need the full functional? . . . . .	11
Remove subscript suppression convention? . . . . .	11
ABOVE NOTATION CREATES AMBIGUITY!!!!!! Check for uses... . . . . .	12
define subset functions, ditch set inputs? appendix? . . . . .	13
italic theta font before Bayes? . . . . .	14
use todonotes package instead of my initials?	
equation numbers to final line!	
brackets for expectation ops?	
line break symbol format, before/after?	
learners (full) or decision functions (range)?	
likelihood function terminology	
Dirichlet localization or concentration?	
fix arguments for theta, nbar dists.	
suppress arguments where sensible? eg $P_d = \theta$	
ALL figure notation: theta font + Ycal indexing. use R,f opt?	
is $\Theta$ redundant given $\mathcal{P}$ ?	
DIM and PR operator from AMS?	
multinomial random process?	
introduce tilde alpha to match tilde theta?	
use x,y subs instead of prime/tilde?	
change from nbar to empirical RV???	
investigate N lim for nbar given theta, bayes risk vs model support	
SEMI-SUPERVISED - generalize to joint decisions!!! training/test!	
priors = sparse conditionals; w/ sufficient statistics	
NFLT investigation? try matlab examples	
use clairvoyant/irreducible terms and symbols	

empirical risk terms/discussion?

check notation throughout for expectation/variance operators, proper arguments

generalize  $y, x, h$  from scalars to functions!!!

jeffrey prior, fisher info?

PGR: bibliography

Theo: (mult moments), DP agg, Dir moments

Bishop: (dir eq), dir posterior, moments, mode

Ferguson: (agg Dir), agg DP - via theo, DP posterior, moments

Gershman: agg DP - ref ferg, discrete draws

Johnson GET PDF: mult moments, (mult agg, DM agg, DM moments)

Add Theo-PR???

# Chapter 1

## Introduction

### 1.1 Background

PGR: complete rework

This report details a Bayesian perspective on statistical learning theory for when both the observations and unobserved quantities are jointly distributed according to an unknown probability distribution function. While the validity of Bayesian methods for statistical signal processing and machine learning has long been contended, the author believes it to be a justified approach that does not necessarily imply that the distribution model is ‘random’; rather, it simply reflects the desire of the user to formulate risk as a weighted sum of learner performance across the space of distributions.

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 informative prior [4] is chosen that is concentrated around the actual data probability distribution, low risk learning functions are possible even with limited training data; however, if the informative prior is poorly designed, a good solution may not be achieved. Conversely, a non-informative prior that weights the different distributions without preference provides a more robust solution for all models, but may underperform relative to learners based on well-selected informative priors.

This work assumes that the prior distribution is Dirichlet. The class of Dirichlet probability density functions (PDF) and processes have the desirable properties of full support over the set of possible data-generating distributions and an analytic posterior distribution for independently and identically distributed data [7]. Furthermore, control of the Dirichlet parameters can enable both non-informative and informative prior knowledge. Special cases including the uniform prior will be given specific attention.

After introducing the problem and discussing the relevant data probability distributions, the Bayesian framework will be applied to two of the most common loss functions in machine learning: the squared error loss function (common for regression), and the 0-1 loss function [1] (common for classification). Optimal estimators/classifiers and their corresponding minimum risk will be presented for different Dirichlet prior distributions. Specific attention will be given to various asymptotic cases to show the differing performance for non-informative and informative Dirichlet priors.

## 1.2 Notation

This section details the mathematical notation and typesetting conventions used throughout. Note that many variable scalars and functions including  $x$ ,  $y$ ,  $g$ , etc. are repeatedly redefined and reused to avoid introducing an excessive volume of symbols; unless explicitly stated, none of these variable definitions will hold in subsequent sections.

### Sets and Function Arguments

Sets will typically be typeset with a calligraphic font, such as  $\mathcal{X}$ . Exceptions include common number sets such as the real numbers, which are typeset using blackboard bold  $\mathbb{R}$ . Function spaces such as the set of functions  $\mathcal{X} \mapsto \mathcal{Y}$  are compactly represented as  $\mathcal{Y}^{\mathcal{X}}$ .

Various mappings will be defined for which the domain and/or the range [16] are function spaces. The set of functions  $\mathcal{X} \mapsto \mathcal{Y}$  is denoted  $\mathcal{Y}^{\mathcal{X}}$ . For a mapping  $g : \mathcal{Z} \mapsto \mathcal{Y}^{\mathcal{X}}$ , the argument notation  $g(z) \in \mathcal{Y}^{\mathcal{X}}$  denotes a function, while  $g(x; z) \in \mathcal{Y}$  is a specific value of that function. Semicolons are used to distinguish between the different groups of arguments. A set of finite sequences  $\{1, \dots, N\} \mapsto \mathcal{S}$  will be represented as  $\mathcal{S}^N$  for brevity.

The convention adopted for natural numbers is  $\mathbb{N} = \{1, 2, \dots\}$ ; the set of non-negative integers is denoted  $\mathbb{Z}_{\geq 0} = \mathbb{N} \cup \{0\}$ . The set of positive real numbers  $\mathbb{R}^+$  excludes zero, while non-negative real numbers are represented as  $\mathbb{R}_{\geq 0} = \mathbb{R}^+ \cup \{0\}$ . The cardinality of countably infinite sets, including the set of natural numbers, is denoted  $\aleph_0 = |\mathbb{N}|$ ; the cardinality of uncountable sets such as  $\mathbb{R}$  is at least  $\aleph_1$ .

Numerous probability distribution functions will be defined over different domains. As such, for a given set  $\mathcal{X}$ , define a set function  $\mathcal{P}$  such that  $\mathcal{P}(\mathcal{X})$  is the set of distributions over  $\mathcal{X}$ . If  $\mathcal{X}$  is countable, the set is defined as  $\mathcal{P}(\mathcal{X}) = \{p \in \mathbb{R}_{\geq 0}^{\mathcal{X}} : \sum_{x \in \mathcal{X}} p(x) = 1\}$ ; if  $\mathcal{X}$  is a Euclidean space, the set is defined as  $\mathcal{P}(\mathcal{X}) = \{p \in \mathbb{R}_{\geq 0}^{\mathcal{X}} : \int_{\mathcal{X}} p(x) dx = 1\}$ .

aleph reference?

## Random elements, variables, and processes

Random elements are denoted with roman font (e.g.  $x$ ), while specific values are denoted with italics (e.g.  $x$ ). Random elements that assume numerical scalars/functions are referred to as random variables/processes, respectively.

Consider a random element  $x \in \mathcal{X}$ . If  $\mathcal{X}$  is countable, either finite with  $|\mathcal{X}| \in \mathbb{N}$  or countably infinite with  $|\mathcal{X}| = \aleph_0$ , then  $x$  is a discrete random element and is characterized by a probability mass function (PMF) [13], denoted  $P_x \in \mathcal{P}(\mathcal{X})$ . If  $\mathcal{X}$  is a Euclidean space and is thus uncountable with  $|\mathcal{X}| \geq \aleph_1$ , then  $x$  is a continuous random variable/process characterized by a probability density function (PDF), denoted  $p_x \in \mathcal{P}(\mathcal{X})$ .

explicit PMF/PDF formula with P of events?

For notational simplicity, probability distributions are occasionally represented as  $P(x)$ ; in such instances, the formal notation can be recovered by replacing the distribution with  $P_x(x)$  and all instances of the roman symbol  $x$  with the italic symbol  $x$ .

Consider  $x$  conditioned on another random element  $z \in \mathcal{Z}$ . The conditional distribution is represented as  $P_{x|z} : \mathcal{Z} \mapsto \mathcal{P}(\mathcal{X})$ , such that  $P_{x|z}(z)$  is a PMF over  $\mathcal{X}$  and  $P_{x|z}(x|z)$  is a specific value of that PMF. These distributions may be compactly represented as  $P(x|z)$ .

Often, the dependency on the conditional variable  $z$  will not be expressed in terms of a specific value  $z$ , but will be left in terms of the random element itself; in this case, the more

compact notation  $P_{x|z}$  is used to imply  $P_{x|z}(z)$ , a function of  $z$ .

Ever need the full functional?

Remove subscript suppression convention?

Many distributions will be repeatedly used and thus special functions will be defined for the PDF's and PMF's of interest. For example, consider a random process  $x \in \mathcal{X}$  characterized by a Dirichlet distribution with parameters  $\alpha \in \mathcal{A}$ ; the PDF will be notated as  $\text{Dir} : \mathcal{A} \mapsto \mathcal{P}(\mathcal{X})$ , where the range is the set of valid PDF's. More compactly, the notation  $x \sim \text{Dir}(\alpha)$  implies that  $P_x = \text{Dir}(\cdot; \alpha)$ . Other distribution functions repeatedly used include Multi, DM, DP, and DMP, representing the multinomial distribution, the Dirichlet-multinomial distribution, the Dirichlet process, and the Dirichlet-Multinomial process.

## Expectation Operators

For a discrete random element  $x$ , the expectation operator  $E_x$  is defined as

$$E_x [g(x)] = \sum_x P_x(x)g(x) , \quad (1.1)$$

where the argument  $g$  is an arbitrary scalar function of  $x$  with range  $\mathbb{R}$ . Additionally, define the variance operator  $C_x$  as

$$C_x [g(x)] = E_x \left[ \left( g(x) - E_x [g(x)] \right)^2 \right] . \quad (1.2)$$

When  $x$  is a random variable and the function  $g$  is the identity operator, such that  $g(x) = x$ , the mean and variance are compactly represented as  $\mu_x$  and  $\Sigma_x$ , respectively.

These operations can be performed with respect to a conditional distribution as well. In this case, the expectation operator is a function of the observed value of  $z$ , such that

$$E_{x|z} [g(x)](z) = \sum_x P_{x|z}(x|z)g(x) . \quad (1.3)$$

Similarly, the conditional variance is notated  $C_{x|z} [g(x)](z)$ . When  $g$  is the identity operator, the conditional mean and variance as represented by  $\mu_{x|z}(z)$  and  $\Sigma_{x|z}(z)$ , respectively.

As with conditional distributions, it is common that an explicit value  $z$  of the conditional random element will not be used, but rather the expectation will be left as a function of the

random element  $z$ . In these cases, the argument is suppressed and the notation  $E_{x|z} [g(x)]$  implies the dependency on  $z$ . This convention also holds for the conditional variance operator  $C_{x|z}$ , as well as for the  $\mu_{x|z}$  and  $\Sigma_{x|z}$  operators.

If the range of  $g$  is a Hilbert space, such that  $g(x)$  is itself a function with a domain  $\mathcal{Y}$ , then the notation for these operators is expanded. The output of the expectation operator is a function over  $\mathcal{Y}$  represented by

$$E_x [g(x)](y) = \sum_x P_x(x) g(y; x) . \quad (1.4)$$

Similarly, the covariance operator notation is modified and the output is a function over  $\mathcal{Y} \times \mathcal{Y}$ ,

$$\begin{aligned} C_x [g(x)](y, y') \\ = E_x \left[ (g(y; x) - E_x [g(y; x)]) (g(y'; x) - E_x [g(y'; x)]) \right] . \end{aligned} \quad (1.5)$$

As before, the notation is simplified when the function  $g$  is the identity operator. If  $x$  is a random process over a domain  $\mathcal{Y}$ , then the mean and covariance functions are  $\mu_x(y)$  and  $\Sigma_x(y, y')$ .

If the expectations are evaluated with respect to a conditional distribution  $P_{x|z}$ , the additional argument for the observed random element is added and the notation for the above operators extends to  $E_{x|z} [g(x)](y|z)$  and  $C_{x|z} [g(x)](y, y'|z)$  for non-scalar outputs. When  $g$  is the identity operator, the notation  $\mu_{x|z}(y|z)$  and  $\Sigma_{x|z}(y, y'|z)$  is used.

Again, it is common for the conditional random element  $z$  to be left as a random quantity instead of being explicitly defined. In such cases, the italic  $z$  is dropped from the arguments and the formulae  $E_{x|z} [g(x)](y)$ ,  $C_{x|z} [g(x)](y, y')$ ,  $\mu_{x|z}(y)$ , and  $\Sigma_{x|z}(y, y')$  imply dependence on  $z$ .

ABOVE NOTATION CREATES AMBIGUITY!!!!!! Check for uses...

As for probability distributions, the subscript notation of these operators may be suppressed. In such cases, the expectations are to be performed with respect to the joint distribution of all random elements (roman font) found in the argument. For example,  $E [f(y, x)| z]$  compactly represents  $E_{y,x|z} [f(y, x)]$ .

## Special Functions

Certain specialized functions are detailed next. Both the Dirac and Kronecker delta functions will be used throughout. The Dirac delta function over a Euclidean domain  $\mathcal{X}$  is represented as  $\delta(\cdot)$ ; it has support only at the point  $x = 0$  and satisfies

$$\int_{\mathcal{X}} \delta(x) dx = 1 . \quad (1.6)$$

Consequently, it also satisfies

$$\int_{\mathcal{X}} g(x) \delta(x) dx = g(0) . \quad (1.7)$$

Consider a countable set  $\mathcal{X}$ ; the Kronecker delta function has domain  $\mathcal{X} \times \mathcal{X}$  and is defined as

$$\delta[x, x'] = \begin{cases} 1 & \text{if } x = x', \\ 0 & \text{if } x \neq x'. \end{cases} \quad (1.8)$$

PGR: reference Dirac/Kronecker?

The multinomial coefficient and multivariate beta function, which typically operate on sequences, are defined more generally for function inputs. The multinomial operator  $\mathcal{M}$  is used for functions  $g : \mathcal{X} \mapsto \mathbb{Z}_{\geq 0}$  that map to non-negative integers from an arbitrary countable domain  $\mathcal{X}$ . The output of the operator is

$$\mathcal{M}(g) = \frac{(\sum_{x \in \mathcal{X}} g(x))!}{\prod_{x \in \mathcal{X}} g(x)!} . \quad (1.9)$$

Similarly, the beta function  $\beta$  operates on functions  $g : \mathcal{X} \mapsto \mathbb{R}^+$  that map to positive real numbers from an arbitrary countable domain  $\mathcal{X}$ , such that

$$\beta(g) = \frac{\prod_{x \in \mathcal{X}} \Gamma(g(x))}{\Gamma(\sum_{x \in \mathcal{X}} g(x))} . \quad (1.10)$$

Note that the countable domains of the input functions may have an infinite number of elements. These functions will also be used to operate on a subset of a functions' domain  $\mathcal{S} \subset \mathcal{X}$  and its corresponding image. Set notation for a function is used to express the argument, so that  $\mathcal{M}\left(\{g(x) : x \in S\}\right)$  and  $\beta\left(\{g(x) : x \in S\}\right)$ .

define subset functions, ditch set inputs? appendix?

# Chapter 2

## Problem Statement

### 2.1 Data Model and Objective

italic theta font before Bayes?

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  $\theta \in \Theta = \mathcal{P}(\mathcal{Y} \times \mathcal{X})$ , such that  $P_{y,x|\theta} = \theta$ . Note that the uppercase PMF notation used throughout this section implies that the random elements are discrete; PDF's are used when  $x$  and/or  $y$  are continuous random variables/processes.

PGR: consider definition/equivalence of D and Y,X. Think indexing.

Also observed is a random sequence of  $N$  samples from  $\theta$ , denoted  $D \in \mathcal{D} = \{\mathcal{Y} \times \mathcal{X}\}^N$ ; an alternative representation that can be used is  $D \equiv (Y, X)$ . The  $N$  data pairs are conditionally independent from one another and are identically distributed as  $P_{D_n|\theta} = \theta$ . The samples are also conditionally independent from  $(y, x)$ . Thus,

$$P_{y,x,D|\theta}(y, x, D|\theta) = P_{y,x|\theta}(y, x|\theta) \prod_{n=1}^N P_{D_n|\theta}(Y_n, X_n|\theta). \quad (2.1)$$

The task in supervised machine learning is to design a decision function  $f : \mathcal{D} \mapsto \mathcal{H}^\mathcal{X}$  which produces a mapping from the space of the observed random elements to a decision space  $\mathcal{H}$ . Define the function space  $\mathcal{F} = \{\mathcal{H}^\mathcal{X}\}^{\mathcal{D}}$ , such that  $f \in \mathcal{F}$ . The learning functions are non-parametric and there are no restrictions on the set of achievable functions  $\mathcal{F}$ .

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 objective is to minimize the conditional expected loss, or conditional “risk”,

$$\begin{aligned}\mathcal{R}_\Theta(f; \theta) &= E_{y,x,D|\theta} \left[ \mathcal{L}(f(x; D), y) \right] \\ &= E_{x,D|\theta} \left[ E_{y|x,\theta} \left[ \mathcal{L}(f(x; D), y) \right] \right] \\ &= E_{D|\theta} \left[ E_{x|\theta} \left[ E_{y|x,\theta} \left[ \mathcal{L}(f(x; D), y) \right] \right] \right].\end{aligned}\tag{2.2}$$

where the conditional independence of random element  $y$  from the training data  $D$  given the model  $\theta$  is used. As the model  $\theta$  is not observed,  $\mathcal{R}_\Theta : \Theta \mapsto \mathbb{R}_{\geq 0}^{\mathcal{F}}$  is not a feasible objective function for optimization. This is the fundamental challenge of supervised learning: the true risk objective is unknown and the designer can never be precisely sure how well any decision function performs.

### 2.1.1 Clairvoyant Decision

PGR: subscript Theta? Use theta sub and remove argument like a cond dist?

PGR: use marginal/conditional thetas?

It is instructive to formulate the optimal decision function assuming the model  $\theta$  was in fact observed; it will be referred to as the “clairvoyant” function, following terminology used in [10]. This clairvoyant decision function  $f_\Theta : \Theta \mapsto \mathcal{F}$  is represented by

$$f_\Theta(\theta) = \arg \min_{f \in \mathcal{F}} \mathcal{R}_\Theta(f; \theta).\tag{2.3}$$

For a given set of observations  $x$  and  $D$ , the function  $f_\Theta(\theta) \in \mathcal{F}$  selects the decision  $h = \arg \min_{h \in \mathcal{H}} E_{y|x,\theta} [\mathcal{L}(h, y)]$ . Note the conditional independence of  $y$  from  $D$  in (2.2) - the knowledge of  $\theta$  renders the training data  $D$  useless. As such, the range of the clairvoyant function is recast as  $f_\Theta : \Theta \mapsto \mathcal{H}^{\mathcal{X}}$  and the decisions are

$$f_\Theta(x; \theta) = \arg \min_{h \in \mathcal{H}} E_{y|x,\theta} [\mathcal{L}(h, y)].\tag{2.4}$$

The corresponding clairvoyant risk for a given model  $\theta$  is

$$\begin{aligned}\mathcal{R}_\Theta^*(\theta) &\equiv \mathcal{R}_\Theta(f_\Theta(\theta); \theta) \\ &= \min_{f \in \mathcal{F}} \mathcal{R}_\Theta(f; \theta) \\ &= E_{x|\theta} \left[ \min_{h \in \mathcal{H}} E_{y|x,\theta} [\mathcal{L}(h, y)] \right] .\end{aligned}\tag{2.5}$$

### 2.1.2 Bayes Decision

To design an optimal decision function  $f \in \mathcal{F}$ , an operator must be chosen to remove the dependency of the conditional risk on  $\theta$  and form an objective function  $\mathcal{F} \mapsto \mathbb{R}_{\geq 0}$ . One choice is to integrate over  $\Theta$ ; to ensure a non-negative objective value, the weighting function should be non-negative. Also, as scaling the objective function will not change its minimizing argument, the weighting function can be constrained to integrate to one. These are the requirements for a valid probability density function (PDF); as such, the model  $\theta$  is treated as a random process and a Bayesian approach can be adopted.

Define the PDF  $p_\theta \in \mathcal{P}(\Theta)$ . Now the Bayes risk can be formulated as

$$\begin{aligned}\mathcal{R}(f) &= E_\theta [\mathcal{R}_\Theta(f; \theta)] \\ &= E_{y,x,D} [\mathcal{L}(f(x; D), y)] \\ &= E_{x,D} \left[ E_{y|x,D} [\mathcal{L}(f(x; D), y)] \right] \\ &= E_D \left[ E_{x|D} \left[ E_{y|x,D} [\mathcal{L}(f(x; D), y)] \right] \right]\end{aligned}\tag{2.6}$$

and  $y$ ,  $x$ , and  $D$  are treated as jointly distributed random elements. Observe that the Bayesian predictive distributions can be represented as  $P_{x|D} = E_{\theta|D} [P_{x|\theta}]$  and  $P_{y|x,D} = E_{\theta|x,D} [P_{y|x,\theta}]$ , the expected values of the corresponding clairvoyant distributions with respect to the model posteriors  $p_{\theta|D}$  and  $p_{\theta|x,D}$ , respectively.

PGR: BELOW, add formula for  $f(D)$ ?

Finally, express the optimal learning function

$$f^* = \arg \min_{f \in \mathcal{F}} \mathcal{R}(f) .\tag{2.7}$$

The decision expressed by the learning function  $f^*$  given observed values of  $x$  and  $D$  is

$$\begin{aligned} f^*(x; D) &= \arg \min_{h \in \mathcal{H}} E_{y|x,D} [\mathcal{L}(h, y)] \\ &= \arg \min_{h \in \mathcal{H}} E_{\theta|x,D} \left[ E_{y|x,\theta} [\mathcal{L}(h, y)] \right]. \end{aligned} \quad (2.8)$$

Thus, the Bayesian approach uses the model posterior  $p_{\theta|x,D}$  to integrate out the dependency on the model given the observable random elements. The minimum Bayes risk is

$$\begin{aligned} \mathcal{R}^* &\equiv \mathcal{R}(f^*) \\ &= \min_{f \in \mathcal{F}} \mathcal{R}(f) \\ &= E_{x,D} \left[ \min_{h \in \mathcal{H}} E_{y|x,D} [\mathcal{L}(h, y)] \right] \\ &= E_D \left[ E_{x|D} \left[ \min_{h \in \mathcal{H}} E_{y|x,D} [\mathcal{L}(h, y)] \right] \right]. \end{aligned} \quad (2.9)$$

### 2.1.2.1 Irreducible Risk

PGR: RECONSIDER IRREDUCIBLE TERMINOLOGY!!!!!!

PGR: bound met in limit of  $N$ ? full support? bounded prior? change/move to clairvoyant discussion?

The clairvoyant risk (2.5) for a given model satisfies  $\mathcal{R}_\Theta(\theta) \leq \mathcal{R}_\Theta(f; \theta) \quad \forall f \in \mathcal{F}, \theta \in \Theta$ . Consequently, the Bayes risk satisfies  $E_\theta [\mathcal{R}_\Theta(\theta)] \leq \mathcal{R}(f) \quad \forall f \in \mathcal{F}$ ; the expected value of the clairvoyant risk will thus be referred to as the “irreducible” risk.

It is important to note that this inequality holds for any number of training samples  $N$  and that the irreducible risk does not depend on  $N$ . Thus, even with unlimited training data, no learning function can provide a Bayes risk lower than this value.

## 2.2 Sufficient Statistic: the Empirical PMF

PGR: MOVE SECTION before Bayesian approach??

PGR: continuous? DMP?

PGR: change to emp PMF RP  $\psi$ ?

PGR: trends in limit of  $N$ , use mean/cov?? EMPIRICAL RISK, bounded prior

For countable sets  $\mathcal{Y}$  and  $\mathcal{X}$ , the distribution of  $D$  conditioned on the model can be formulated as

$$\begin{aligned} P_{D|\theta}(D|\theta) &= \prod_{n=1}^N P_{D_n|\theta}(D_n|\theta) \\ &= \prod_{y \in \mathcal{Y}} \prod_{x \in \mathcal{X}} \theta(y, x)^{\bar{N}(y, x; D)}, \end{aligned} \quad (2.10)$$

where the dependency on the training data  $D$  is expressed through a transform function  $\bar{N} : \mathcal{D} \mapsto \bar{\mathcal{N}}$ , where the range is

$$\bar{\mathcal{N}} = \left\{ \bar{n} \in \mathbb{Z}_{\geq 0}^{\mathcal{Y} \times \mathcal{X}} : \sum_{y \in \mathcal{Y}} \sum_{x \in \mathcal{X}} \bar{n}(y, x) = N \right\} \quad (2.11)$$

and the function is defined as

$$\begin{aligned} \bar{N}(y, x; D) &= \sum_{n=1}^N \delta[(y, x), D_n] \\ &= \sum_{n=1}^N \delta[y, Y_n] \delta[x, X_n]. \end{aligned} \quad (2.12)$$

This function counts the number of occurrences of the pair  $(y, x)$  in the training set  $D$ .

PGR: show SS via Kay, data likelihood??!!

Note that  $P_{D|\theta}$  depends on the training data  $D$  only through the transform  $\bar{N}$ ;  $\bar{N}(D)$  is thus a sufficient statistic [2] for the model  $\theta$ . Consequently, other distributions of interest  $P_D$ ,  $P_{x|D}$ , and  $P_{y|x,D}$  will also depend on  $D$  via  $\bar{N}(D)$ . As such, it is useful to define a new random process  $\bar{n} \equiv \bar{N}(D) \in \bar{\mathcal{N}}$ .

Frequently, the corresponding distributions  $P_{\bar{n}}$ ,  $P_{x|\bar{n}}$ , and  $P_{y|x,\bar{n}}$  will be used to find the optimal decision functions and the minimum risk. Note that  $\mathcal{M}(\bar{N}(D)) P_{D|\theta}(D|\theta) = P_{\bar{n}|\theta}(\bar{N}(D)|\theta)$ , where  $\mathcal{M}$  is the multinomial operator. Also note that  $P_{x|D}(D) = P_{x|\bar{n}}(\bar{N}(D))$  and  $P_{y|x,D}(x, D) = P_{y|x,\bar{n}}(x, \bar{N}(D))$ .

The cardinality of the random process' domain is  $|\bar{\mathcal{N}}| = \mathcal{M}(\{N, |\mathcal{Y}| |\mathcal{X}| - 1\})$ ; this can be shown using the stars-and-bars method [6]. 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. Also, observe that the set  $\{\bar{n}/N : \bar{n} \in \bar{\mathcal{N}}\} \subset \Theta$  and thus that the empirical distribution  $\bar{N}(D)/N$  assumes one of a finite number of the elements from  $\Theta$ .

PGR: sufficient statistic savings in memory bits???

Conditioned on the model  $\theta$ , the PMF of  $\bar{n}$  is a multinomial distribution

$$\begin{aligned} P_{\bar{n}|\theta}(\bar{n}|\theta) &= \sum_{D:\bar{N}(D)=\bar{n}} P_{D|\theta}(D|\theta) \\ &= |\{D : \bar{N}(D) = \bar{n}\}| \prod_{y \in \mathcal{Y}} \prod_{x \in \mathcal{X}} \theta(y, x)^{\bar{n}(y,x)} \\ &= \mathcal{M}(\bar{n}) \prod_{y \in \mathcal{Y}} \prod_{x \in \mathcal{X}} \theta(y, x)^{\bar{n}(y,x)} \\ &= \text{Multi}(\bar{n}; N, \theta), \end{aligned} \quad (2.13)$$

where the multinomial operator  $\mathcal{M}$  is used.

The first and second joint moments of this multinomial distribution are [18]

$$\mu_{\bar{n}|\theta} = N\theta \quad (2.14)$$

and

$$\begin{aligned} E_{\bar{n}|\theta} [\bar{n}(y, x)\bar{n}(y', x')] &= N(\theta(y, x)\delta[y, y']\delta[x, x'] + (N-1)\theta(y, x)\theta(y', x')) \\ &= N(\theta(y, x)\delta[y, y']\delta[x, x'] - \theta(y, x)\theta(y', x')). \end{aligned} \quad (2.15)$$

and the covariance function is

$$\Sigma_{\bar{n}|\theta}(y, x, y', x') = N(\theta(y, x)\delta[y, y']\delta[x, x'] - \theta(y, x)\theta(y', x')). \quad (2.16)$$

PGR: trends with N? delta? Figures??

Also, observe that the maximum likelihood estimate of  $\theta$  given the training statistic is [14],

$$\begin{aligned} \theta_{ML}(\bar{n}) &= \arg \max_{\theta \in \Theta} P_{\bar{n}|\theta}(\bar{n}|\theta) \\ &= \frac{\bar{n}}{N}, \end{aligned} \quad (2.17)$$

the empirical distribution.

## 2.3 Marginal/Conditional Distributions for Observed/Unobserved Elements

PGR: move section?

PGR: clean up

### 2.3.1 Marginal and Conditional Distributions of $\theta$

PGR: are they independent??? No?!

PGR: bijection!!

As only  $y$  is unobservable, it will be useful to decompose the model distribution as  $\theta \equiv (\theta', \tilde{\theta})$ . First, introduce the marginal distribution  $\theta' \equiv \sum_{y \in \mathcal{Y}} \theta(y, \cdot) \in \mathcal{P}(\mathcal{X})$ ; note that the summation is replaced by an integral when  $y$  is a continuous random variable. Next, introduce the conditional distributions  $\tilde{\theta} \in \mathcal{P}(\mathcal{Y}|\mathcal{X})$  defined as  $\tilde{\theta}(x) \equiv \theta(\cdot, x)/\theta'(x)$ .

This decomposition enables the clairvoyant distributions to be represented as  $P_{x|\theta} = P_{x|\theta'} = \theta'$  and  $P_{y|x,\theta} = P_{y|x,\tilde{\theta}} = \tilde{\theta}(x)$ ; these distributions will be of recurring importance.

PGR: conditional theta condition on  $x$  necessary above? Yes?!

PGR: marginal theta conditional on  $X$ , not full  $D$  above???? No!?

### 2.3.2 Marginal and Conditional Distributions of Training Data

Also of interest are the marginal and conditional distributions of the joint training data partitions  $Y$  and  $X$ . The marginal distribution given  $\theta$  for the observations  $X$  alone is

$$\begin{aligned} P_{X|\theta}(X|\theta) &= \prod_{n=1}^N P_{X_n|\theta}(X_n|\theta) \\ &= \prod_{x \in \mathcal{X}} \theta'(x)^{N'(x; X)}, \end{aligned} \tag{2.18}$$

where the dependency on  $\theta$  is only through the marginal model  $\theta'$ . Additionally, note that the dependency on the training observations  $X$  is expressed through a “marginal” counting function  $N' : \mathcal{X} \mapsto \mathcal{N}'$  with range

$$\mathcal{N}' = \left\{ n' \in \mathbb{Z}_{\geq 0}^{\mathcal{X}} : \sum_{x \in \mathcal{X}} n'(x) = N \right\}, \tag{2.19}$$

defined as

$$N'(X) = \sum_{n=1}^N \delta[\cdot, X_n] \equiv \sum_{y \in \mathcal{Y}} \bar{N}(y, \cdot; D). \tag{2.20}$$

The conditional distribution of the values  $Y$  given the corresponding  $X$  can be found using Bayes theorem as

$$\begin{aligned} P_{Y|X,\theta}(Y|X,\theta) &= \prod_{n=1}^N \frac{P_{Y_n,X_n|\theta}(Y_n, X_n|\theta)}{P_{X_n|\theta}(X_n|\theta)} \\ &= \prod_{x \in \mathcal{X}} \left[ \prod_{y \in \mathcal{Y}} \tilde{\theta}(y; x)^{\bar{N}(y,x; Y, X)} \right], \end{aligned} \quad (2.21)$$

which is dependent on the model  $\theta$  only through the conditional models  $\tilde{\theta}(x)$ .

As before, the dependency on the training data can be simplified using a sufficient statistic. Introduce the “marginalized” random process  $n'$  over the set  $\mathcal{X}$ , defined as  $n' \equiv \sum_{y \in \mathcal{Y}} \bar{n}(y, \cdot) \equiv N'(X) \in \mathcal{N}'$ . By the aggregation property of Multinomial random processes [9], the aggregation conditioned on the model  $\theta$  is distributed as  $n'|\theta \sim \text{Multi}(N, \theta')$ .

Also of interest is the distribution of  $\bar{n}$  conditioned on its aggregation  $n'$ . Using the multinomial distribution properties proven in Appendix A.2, it can be shown that when conditioned on the model  $\theta$  as well, the PMF of  $\bar{n}$  is

$$\begin{aligned} P_{\bar{n}|n',\theta}(\bar{n}|n',\theta) &= \prod_{x \in \mathcal{X}} \left[ \mathcal{M}(\bar{n}(\cdot, x)) \prod_{y \in \mathcal{Y}} \tilde{\theta}(y; x)^{\bar{n}(y,x)} \right] \\ &= \prod_{x \in \mathcal{X}} \text{Multi}\left(\bar{n}(\cdot, x); n'(x), \tilde{\theta}(x)\right), \end{aligned} \quad (2.22)$$

over the domain  $\left\{ \bar{n} \in \mathbb{Z}_{\geq 0}^{|\mathcal{Y}| \times |\mathcal{X}| : \sum_{y \in \mathcal{Y}} \bar{n}(y, \cdot) = n'} \right\}$ . Observe that conditioning on the aggregation renders the function segments  $\bar{n}(\cdot, x)$  independent of one another and that they are also Multinomial, such that  $\bar{n}(\cdot, x)|n'(x), \theta \sim \text{Multi}(n'(x), \tilde{\theta}(x))$ . Furthermore, the dependency on  $\theta$  is expressed through the conditional model  $\tilde{\theta}$ .

### 2.3.3 Model Posteriors

PGR: clean, discuss

The Bayesian predictive distributions analogous to the clairvoyant distributions can now be simplified as  $P_{x|D} = \mu_{\theta'|D}$  and  $P_{y|x,D} = \mu_{\tilde{\theta}(x)|x,D}$ .

PGR: use  $P_{y|x,D} = E_{\tilde{\theta}|x,D} [\tilde{\theta}(x)]$  instead???

$$\begin{aligned}
p_{\theta', \tilde{\theta}|Y, X}(\theta', \tilde{\theta}|Y, X) &= \frac{P_{Y|X, \tilde{\theta}}(Y|X, \tilde{\theta})}{P_{Y|X}(Y|X)} \frac{P_{X|\theta'}(X|\theta')}{P_X(X)} p_{\theta', \tilde{\theta}}(\theta', \tilde{\theta}) \\
&= p_{\tilde{\theta}|Y, X}(\tilde{\theta}|Y, X) p_{\theta'|X}(\theta'|X) \frac{p_{\tilde{\theta}|\theta'}(\tilde{\theta}|\theta')}{p_{\tilde{\theta}|X}(\tilde{\theta}|X)}
\end{aligned} \tag{2.23}$$

$$\begin{aligned}
p_{\theta', \tilde{\theta}|\bar{n}}(\theta', \tilde{\theta}|\bar{n}) &= \frac{P_{\bar{n}|n', \tilde{\theta}}(\bar{n}|\sum_y \bar{n}(y, \cdot), \tilde{\theta})}{P_{\bar{n}|n'}(\bar{n}|\sum_y \bar{n}(y, \cdot))} \frac{P_{n'|\theta'}(\sum_y \bar{n}(y, \cdot)|\theta')}{P_{n'}(\sum_y \bar{n}(y, \cdot))} p_{\theta', \tilde{\theta}}(\theta', \tilde{\theta}) \\
&= p_{\tilde{\theta}|\bar{n}, n'}(\tilde{\theta}|\bar{n}, \sum_y \bar{n}(y, \cdot)) p_{\theta'|n'}(\theta'|\sum_y \bar{n}(y, \cdot)) \frac{p_{\tilde{\theta}|\theta'}(\tilde{\theta}|\theta')}{p_{\tilde{\theta}|n'}(\tilde{\theta}|\sum_y \bar{n}(y, \cdot))}
\end{aligned} \tag{2.24}$$

Note that  $p_{\tilde{\theta}|X} = E_{\theta'|X} [p_{\tilde{\theta}|\theta'}]$  and  $p_{\tilde{\theta}|n'} = E_{\theta'|n'} [p_{\tilde{\theta}|\theta'}]$ . Also, note that  $P_{Y|X} = E_{\theta'|X} [E_{\tilde{\theta}|\theta'} [P_{Y|X, \tilde{\theta}}]]$  and  $P_{\bar{n}|n'} = E_{\theta'|n'} [E_{\tilde{\theta}|\theta'} [P_{\bar{n}|n', \tilde{\theta}}]]$ .

$$\begin{aligned}
p_{\theta', \tilde{\theta}|Y, X, x}(\theta', \tilde{\theta}|Y, X, x) &= \frac{P_{Y|X, \tilde{\theta}}(Y|X, \tilde{\theta})}{P_{Y|X, x}(Y|X, x)} \frac{P_{X,x|\theta'}(X, x|\theta')}{P_{X,x}(X, x)} p_{\theta', \tilde{\theta}}(\theta', \tilde{\theta}) \\
&= p_{\tilde{\theta}|Y, X}(\tilde{\theta}|Y, X) p_{\theta'|X, x}(\theta'|X, x) \frac{P_{Y|X}(Y|X)}{P_{Y|X, x}(Y|X, x)} \frac{p_{\tilde{\theta}|\theta'}(\tilde{\theta}|\theta')}{p_{\tilde{\theta}|X}(\tilde{\theta}|X)}
\end{aligned} \tag{2.25}$$

$$\begin{aligned}
p_{\theta', \tilde{\theta}|\bar{n}, x}(\theta', \tilde{\theta}|\bar{n}, x) &= \frac{P_{\bar{n}|n', \tilde{\theta}}(\bar{n}|\sum_y \bar{n}(y, \cdot), \tilde{\theta})}{P_{\bar{n}|n', x}(\bar{n}|\sum_y \bar{n}(y, \cdot), x)} \frac{P_{n', x|\theta'}(\sum_y \bar{n}(y, \cdot), x|\theta')}{P_{n', x}(\sum_y \bar{n}(y, \cdot), x)} p_{\theta', \tilde{\theta}}(\theta', \tilde{\theta}) \\
&= p_{\tilde{\theta}|\bar{n}, n'}(\tilde{\theta}|\bar{n}, \sum_y \bar{n}(y, \cdot)) p_{\theta'|n', x}(\theta'|\sum_y \bar{n}(y, \cdot), x) \\
&\quad \frac{P_{\bar{n}|n'}(\bar{n}|\sum_y \bar{n}(y, \cdot))}{P_{\bar{n}|n', x}(\bar{n}|\sum_y \bar{n}(y, \cdot), x)} \frac{p_{\tilde{\theta}|\theta'}(\tilde{\theta}|\theta')}{p_{\tilde{\theta}|n'}(\tilde{\theta}|\sum_y \bar{n}(y, \cdot))}
\end{aligned} \tag{2.26}$$

Note that  $P_{Y|X, x} = E_{\theta'|X, x} [E_{\tilde{\theta}|\theta'} [P_{Y|X, \tilde{\theta}}]]$  and also that  $P_{\bar{n}|n', x} = E_{\theta'|n', x} [E_{\tilde{\theta}|\theta'} [P_{\bar{n}|n', \tilde{\theta}}]]$ .

## 2.4 Applications to Common Loss Functions

PGR: equation label conflicts??

PGR: marginal/conditional???

PGR: sum vs int???

PGR: move before nbar, keep D?

In this section, loss functions typical for classification and regression applications, specifically the 0-1 loss function and the squared-error loss function, are adopted. The conditional risk (2.2) is assessed, clairvoyant learners (2.4) are found, and the clairvoyant risk (2.5) is expressed.

### 2.4.1 Regression: 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 a closed-form expression of the minimizing estimation function.

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

The loss function is defined as

$$\mathcal{L}(h, y) = (h - y)^2. \quad (2.27)$$

Substituting the squared-error loss into (2.2), the conditional squared-error risk is

$$\begin{aligned} \mathcal{R}_\theta(f; \theta) &= E_{D|\theta} \left[ E_{y,x|\theta} \left[ (f(x; D) - y)^2 \right] \right] \\ &= E_{x|\theta} \left[ E_{y|x,\theta} \left[ E_{D|\theta} \left[ (f(x; D) - y)^2 \right] \right] \right] \\ &= E_{x|\theta} \left[ E_{y|x,\theta} \left[ (y - \mu_{y|x,\theta})^2 \right] \right] + E_{x,D|\theta} \left[ (f(x; D) - \mu_{y|x,\theta})^2 \right] \\ &= E_{x|\theta} \left[ \Sigma_{y|x,\theta} \right] + E_{x,D|\theta} \left[ (f(x; D) - \mu_{y|x,\theta})^2 \right], \end{aligned} \quad (2.28)$$

a sum of two terms. The first term is the expected conditional variance of the true predictive distribution  $P_{y|x,\theta}$ . The second term is the expected squared bias between the Bayesian estimate and the true mean  $\mu_{y|x,\theta}$ .

#### 2.4.1.1 Clairvoyant Estimation

To find the clairvoyant estimator, the squared-error loss is substituted into (2.4); note that the objective function is quadratic over the argument  $h \in \mathcal{H} = \mathbb{R}$ . 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 clairvoyant estimate is the expected value of  $y$  given the model  $\theta$  and the observed value  $x$ , such that

$$\begin{aligned} f_\Theta(x; \theta) &= \arg \min_{h \in \mathbb{R}} E_{y|x,\theta} [(h - y)^2] \\ &= \mu_{y|x,\theta} = \sum_{y \in \mathcal{Y}} y \tilde{\theta}(y; x) . \end{aligned} \quad (2.29)$$

Substituting the loss and clairvoyant function into (2.5), the resulting clairvoyant risk is

$$\begin{aligned} \mathcal{R}_\Theta^*(\theta) &= E_{x|\theta} \left[ E_{y|x,\theta} [(y - \mu_{y|x,\theta})^2] \right] \\ &= E_{x|\theta} [\Sigma_{y|x,\theta}] . \end{aligned} \quad (2.30)$$

Observe that the general conditional risk (2.28) can be represented as  $\mathcal{R}_\Theta(f; \theta) = \mathcal{R}_\Theta^*(\theta) + E_{x,D|\theta} \left[ (f(x; D) - f_\Theta(x; \theta))^2 \right]$ . The first summand is equal to the clairvoyant squared-error; the second term is dependent on the difference between the general estimate and the clairvoyant estimate.

Figure 2.1 displays the clairvoyant risk for predictive models  $\tilde{\theta}(x)$  independent of  $x$ .

#### 2.4.1.2 Bayesian Estimation

##### Optimal Estimate: the Posterior Mean PGR: plots?

To find the optimal estimator, the squared-error loss is substituted into (2.8). Again, 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

$$\begin{aligned} f^*(x; D) &= \arg \min_{h \in \mathbb{R}} E_{y|x,D} [(h - y)^2] \\ &= \mu_{y|x,D} = E_{\theta|x,D} [\mu_{y|x,\theta}] . \end{aligned} \quad (2.31)$$

An interesting form for the optimal estimator is  $f^*(x; D) = E_{\theta|x,D} [f_\Theta(x; \theta)]$ . Substituting the squared-error loss into the second line of (2.8), the optimal Bayes estimator is the conditional expected value of the clairvoyant estimate with respect to the model posterior distribution.

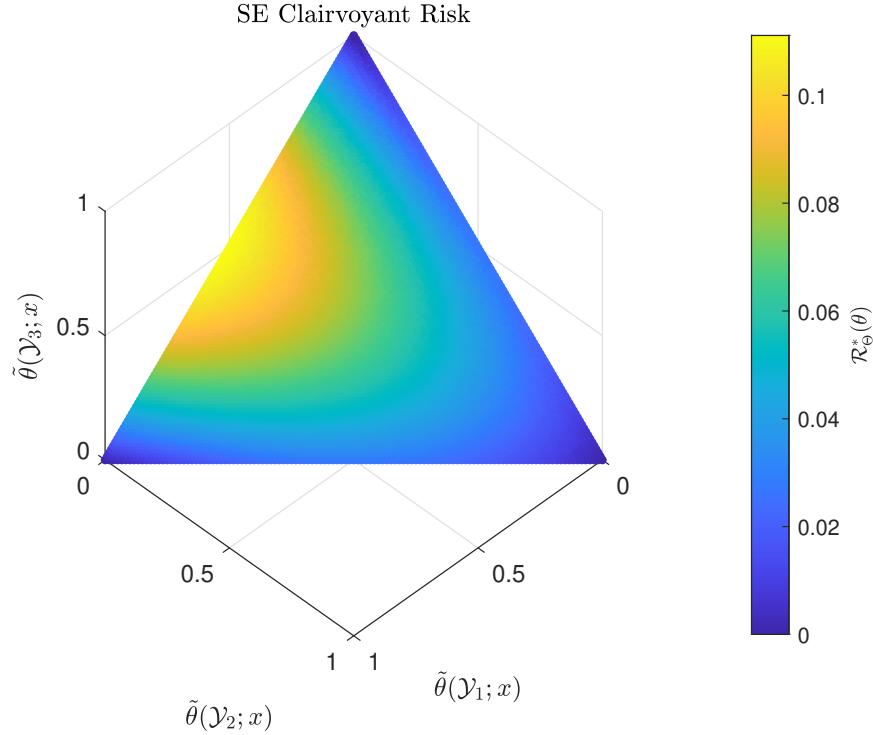


Figure 2.1: Clairvoyant Risk for the Squared-Error Loss Function, constant  $\tilde{\theta}(x)$

**Minimum Risk: the Expected Posterior Variance** The Bayes squared-error risk for a general learning function is

$$\begin{aligned}
 \mathcal{R}(f) &= E_\theta \left[ E_{D|\theta} \left[ E_{y,x|\theta} \left[ (f(x; D) - y)^2 \right] \right] \right] \\
 &= E_{x,D} \left[ E_{y|x,D} \left[ (f(x; D) - y)^2 \right] \right] \\
 &= E_\theta \left[ \mathcal{R}_\Theta^*(\theta) \right] + E_{x,D,\theta} \left[ (f(x; D) - \mu_{y|x,\theta})^2 \right] \\
 &= E_{x,D} \left[ \Sigma_{y|x,D} \right] + E_{x,D} \left[ (f(x; D) - \mu_{y|x,D})^2 \right].
 \end{aligned} \tag{2.32}$$

Substituting the optimal estimator (2.31) into Equation (2.32), the minimum Bayes risk is the expected conditional variance

$$\begin{aligned}
 \mathcal{R}^* &= E_{x,D} \left[ \Sigma_{y|x,D} \right] \\
 &= E_{x,\theta} \left[ \Sigma_{y|x,\theta} \right] + E_{x,D} \left[ C_{\theta|x,D} \left[ \mu_{y|x,\theta} \right] \right] \\
 &= E_\theta \left[ \mathcal{R}_\Theta^*(\theta) \right] + E_{x,D} \left[ C_{\theta|x,D} \left[ f_\Theta(x; \theta) \right] \right].
 \end{aligned} \tag{2.33}$$

The first term is the irreducible risk. The second term is the expected variance of the

clairvoyant estimate  $f_\Theta(x; \theta) = \mu_{y|x,\theta}$  with respect to the model posterior PDF  $p_{\theta|x,D}$

PGR: assess irreducible risk by performing theta agg conditioning???

## 2.4.2 Classification: the 0-1 Loss

In this section, the developed framework is applied to a common machine learning task: classification. In classification problems, the set  $\mathcal{Y}$  is countable and typically finite. Furthermore, the hypothesis space is usually identical to the unobserved variable space; that is  $\mathcal{H} = \mathcal{Y}$ . The 0-1 loss function is the most widely used for these problems; it is represented as

$$\mathcal{L}(h, y) = 1 - \delta[h, y] . \quad (2.34)$$

Applying the 0-1 loss, the conditional risk (2.2) for a general classifier is

$$\begin{aligned} \mathcal{R}_\Theta(f; \theta) &= 1 - E_{D|\theta} \left[ E_{y,x|\theta} \left[ \delta[f(x; D), y] \right] \right] \\ &= 1 - \sum_{x \in \mathcal{X}} E_{D|\theta} \left[ \theta(f(x; D), x) \right] \\ &= 1 - E_{x|\theta} \left[ E_{D|\theta} \left[ P_{y|x,\theta} (f(x; D) | x, \theta) \right] \right] \\ &= 1 - \sum_{x \in \mathcal{X}} \theta'(x) E_{D|\theta} \left[ \tilde{\theta}(f(x; D); x) \right] . \end{aligned} \quad (2.35)$$

### 2.4.2.1 Clairvoyant Hypothesis

To find the clairvoyant classifier, the 0-1 loss is substituted into (2.4); given an observation  $x$ , the optimum hypothesis is simply the value  $y$  that maximizes the conditional model  $\tilde{\theta}(x)$ ,

$$\begin{aligned} f_\Theta(x; \theta) &= \arg \min_{h \in \mathcal{Y}} E_{y|x,\theta} [1 - \delta[h, y]] \\ &= \arg \max_{h \in \mathcal{Y}} P_{y|x,\theta}(h | x, \theta) \\ &= \arg \max_{y \in \mathcal{Y}} \tilde{\theta}(y; x) \\ &= \arg \max_{y \in \mathcal{Y}} \theta(y, x) . \end{aligned} \quad (2.36)$$

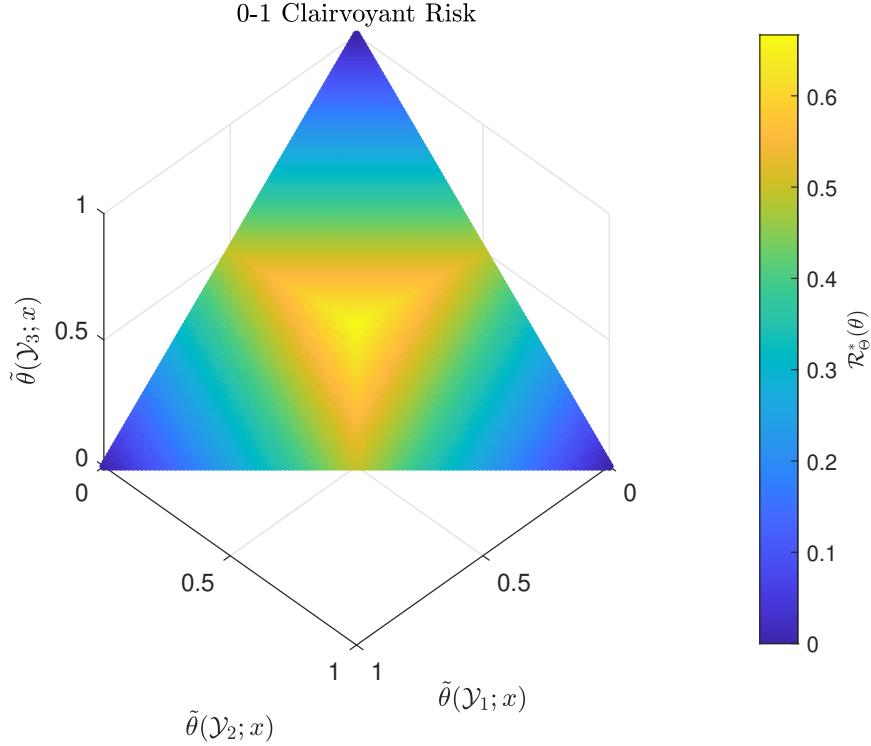


Figure 2.2: Clairvoyant Risk for the 0–1 Loss Function, constant  $\tilde{\theta}(x)$

Substituting the 0-1 loss and clairvoyant hypothesis into (2.5), the resulting clairvoyant risk is

$$\begin{aligned}
 \mathcal{R}_{\Theta}^*(\theta) &= 1 - E_{x|\theta} \left[ \max_{y \in \mathcal{Y}} P_{y|x,\theta}(y|x, \theta) \right] \\
 &= 1 - \sum_{x \in \mathcal{X}} \theta'(x) \max_{y \in \mathcal{Y}} \tilde{\theta}(y; x) \\
 &= 1 - \sum_{x \in \mathcal{X}} \max_{y \in \mathcal{Y}} \theta(y, x) .
 \end{aligned} \tag{2.37}$$

Figure 2.2 displays the clairvoyant risk for predictive models  $\tilde{\theta}(x)$  independent of  $x$ . Intuitively, the models that are more concentrated lead to lower probability of error.

#### 2.4.2.2 Bayesian Classification

**Optimal Hypothesis: Conditional Maximum *a posteriori*** PGR: decision region figures??

To determine the optimal learning function, the 0-1 loss from Equation (2.34) is substi-

tuted into Equation (2.8) to find

$$\begin{aligned} f^*(x; D) &= \arg \min_{h \in \mathcal{Y}} E_{y|x,D} [1 - \delta[h, y]] \\ &= \arg \max_{y \in \mathcal{Y}} P_{y|x,D}(y|x, D). \end{aligned} \quad (2.38)$$

The optimal classifier chooses the value  $y \in \mathcal{Y}$  that maximizes the conditional PMF for the observed values of  $x$  and  $D$ .

**Minimum Risk: Probability of Error** Using the 0-1 loss, the Bayes probability of error (2.9) is

$$\mathcal{R}(f) = 1 - E_{x,D} [P_{y|x,D}(f(x; D)|x, D)]. \quad (2.39)$$

Substituting the optimal learner (2.38) into the general risk (2.39), the minimum probability of error is

$$\mathcal{R}^* = 1 - E_{x,D} \left[ \max_{y \in \mathcal{Y}} P_{y|x,D}(y|x, D) \right]. \quad (2.40)$$

# Chapter 3

## Finite Dirichlet Model

This chapter demonstrates the optimal decision functions when the sets  $\mathcal{Y}$  and  $\mathcal{X}$  have a finite number of elements and the model  $\theta$  is characterized by a Dirichlet distribution.

### 3.1 Probability Distributions

To determine the optimal decision function, the joint PMF  $P_{y,x,D}$  is required. Having already defined the distribution conditioned on the model  $\theta$ , all that remains is to select a PDF  $p_\theta$  reflecting the user's prior knowledge. In this section, the Dirichlet distribution is used. The Dirichlet distribution possesses the desirable property of being the conjugate prior for the multinomial conditional distribution characterizing the data; as such, it will provide analytic forms for the model posterior distribution and lead to closed form expressions for the data conditional distribution used to design the decision function.

Other distributions of interest will be provided, such as the training data PMF  $P_D$  and the conditional distribution  $P_{y|x,D}$  used to form a decision given specific observations.

#### 3.1.1 Model PDF, $p_\theta$

The Dirichlet PDF for the model random process  $\theta \in \Theta$  is [3]

$$\begin{aligned} p_\theta(\theta) &= \beta(\alpha)^{-1} \prod_{y \in \mathcal{Y}} \prod_{x \in \mathcal{X}} \theta(y, x)^{\alpha(y,x)-1} \\ &= \text{Dir}(\theta; \alpha), \end{aligned} \tag{3.1}$$

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

The parameter  $\alpha$  controls around which models  $\theta$  the PDF concentrates and how strongly. For convenience, introduce the concentration parameter  $\alpha_0 \equiv \sum_{y \in \mathcal{Y}} \sum_{x \in \mathcal{X}} \alpha(y, x)$ .

The first and second joint moments of the model are

$$\mu_\theta = \frac{\alpha}{\alpha_0} \quad (3.2)$$

and

$$E_\theta [\theta(y, x)\theta(y', x')] = \frac{\alpha(y, x)\alpha(y', x') + \alpha(y, x)\delta[y, y']\delta[x, x']}{\alpha_0(\alpha_0 + 1)}. \quad (3.3)$$

Observe that  $P_{y,x} = \mu_\theta = \alpha/\alpha_0$ . The covariance is

$$\begin{aligned} \Sigma_\theta(y, x, y', x') &= E_\theta [(\theta(y, x) - \mu_\theta)(\theta(y', x') - \mu_\theta)] \\ &= \frac{\mu_\theta(y, x)\delta[y, y']\delta[x, x'] - \mu_\theta(y, x)\mu_\theta(y', x')}{\alpha_0 + 1}. \end{aligned} \quad (3.4)$$

Also, for  $\alpha(y, x) > 1$ , the maximizing value of the distribution is

$$\theta_{\max} = \arg \max_{\theta \in \Theta} p_\theta(\theta) = \frac{\alpha - 1}{\alpha_0 - |\mathcal{Y}||\mathcal{X}|}. \quad (3.5)$$

This can be easily shown by maximizing the logarithm of the distribution using the method of Lagrange multipliers, as demonstrated in B.1.

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

$$p_\theta(\theta) \rightarrow \delta\left(\theta - \frac{\alpha}{\alpha_0}\right). \quad (3.6)$$

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

$$p_\theta(\theta) \rightarrow \sum_{y \in \mathcal{Y}} \sum_{x \in \mathcal{X}} \frac{\alpha(y, x)}{\alpha_0} \delta(\theta - \delta[\cdot, y]\delta[\cdot, x]), \quad (3.7)$$

which distributes its weight among the  $|\mathcal{Y}||\mathcal{X}|$  models with an  $\ell_0$  norm satisfying  $\|\theta\|_0 = 1$ . Note that the Dirac delta for these formulas is defined on the set  $\Theta$ , such that  $\int_\Theta \delta(\theta)d\theta = 1$ .

PGR: formal proof for limiting PDFs??? stirling/gautschi?

These trends are demonstrated with Figure 3.1. The cardinalities  $|\mathcal{Y}| = 3$  and  $|\mathcal{X}| = 1$  are chosen to enable visualization, despite the implication that  $x$  is deterministic; these cardinalities will be used for many subsequent figures as well. Note that for  $\alpha_0 = 2.99$ ,  $\alpha < 1$  and the PDF values at the boundaries of the domain tend to infinity; this is not captured by the plot color scale.

**Uniform Prior** When the parameterizing function is  $\alpha(y, x) = 1$ , the distribution becomes a uniform PDF and is represented as

$$p_\theta = (|\mathcal{Y}||\mathcal{X}| - 1)! . \quad (3.8)$$

Note that the concentration parameter is  $\alpha_0 = |\mathcal{Y}||\mathcal{X}|$  and  $P_{y,x} = (|\mathcal{Y}||\mathcal{X}|)^{-1}$  is also uniform. Figure 3.2 shows the uniform distribution amplitude for  $|\mathcal{Y}| = 3$  and  $|\mathcal{X}| = 1$ .

### 3.1.1.1 Marginal and Conditional Distributions

PGR: move/add Dir figs here?

The marginal distribution  $\theta'$  and the conditional distribution  $\tilde{\theta}$  will also be of interest. By the aggregation property [7],  $\theta'$  is a Dirichlet random process parameterized by  $\alpha' : \mathcal{X} \mapsto \mathbb{R}^+$ , where  $\alpha' \equiv \sum_{y \in \mathcal{Y}} \alpha(y, \cdot)$ . Note that  $P_x = \mu_{\theta'} = \alpha'/\alpha_0$ .

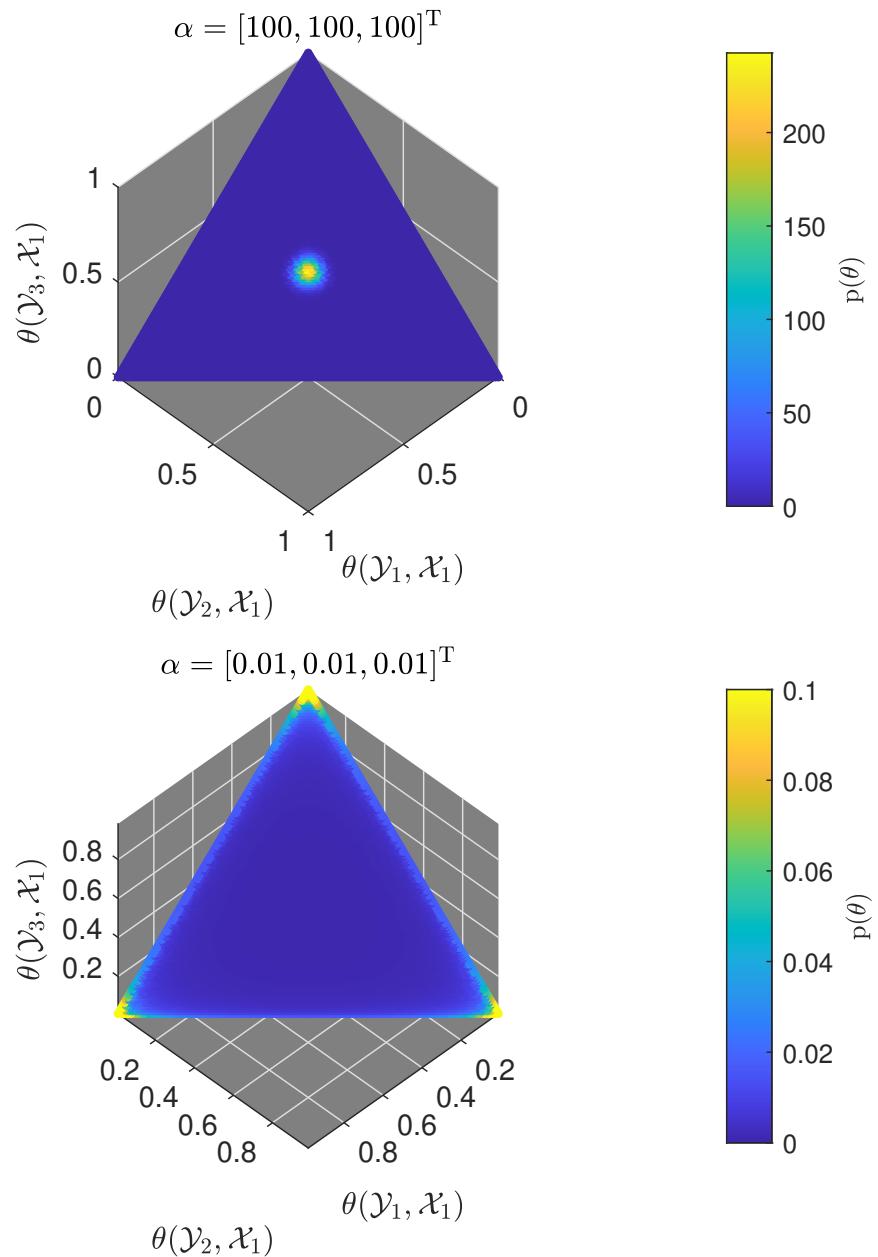
PGR: introduce tilde alpha to match tilde theta?

Also of interest is the distribution of the predictive model  $\tilde{\theta}$  conditioned on the marginal  $\theta'$ . As demonstrated in Appendix A.1, these random processes are jointly distributed as

$$\begin{aligned} p_{\tilde{\theta}|\theta'}(\tilde{\theta}|\theta') &= \prod_{x \in \mathcal{X}} \left[ \beta(\alpha(\cdot, x))^{-1} \prod_{y \in \mathcal{Y}} \tilde{\theta}(y; x)^{\alpha(y, x)-1} \right] \\ &= \prod_{x \in \mathcal{X}} \text{Dir}(\tilde{\theta}(x); \alpha(\cdot, x)) , \end{aligned} \quad (3.9)$$

a product of Dirichlet distributions defined on  $\tilde{\theta} \in \mathcal{P}(\mathcal{Y})^{\mathcal{X}}$ . As shown, the processes  $\tilde{\theta}(x)$  are Dirichlet with parameterizing functions  $\alpha(\cdot, x)$ , independent of one another, and independent of the marginal distribution  $\theta'$ . Observe that the values  $\alpha'(x)$  represent the concentration parameters for the individual Dirichlet processes; also, note that  $P_{y|x} = \mu_{\tilde{\theta}(x)} = \alpha(\cdot, x)/\alpha'(x)$ .

PGR: use conditional independence property to simplify throughout?!? In loss app sections?!?

Figure 3.1: Model prior PDF for different concentrations  $\alpha_0$

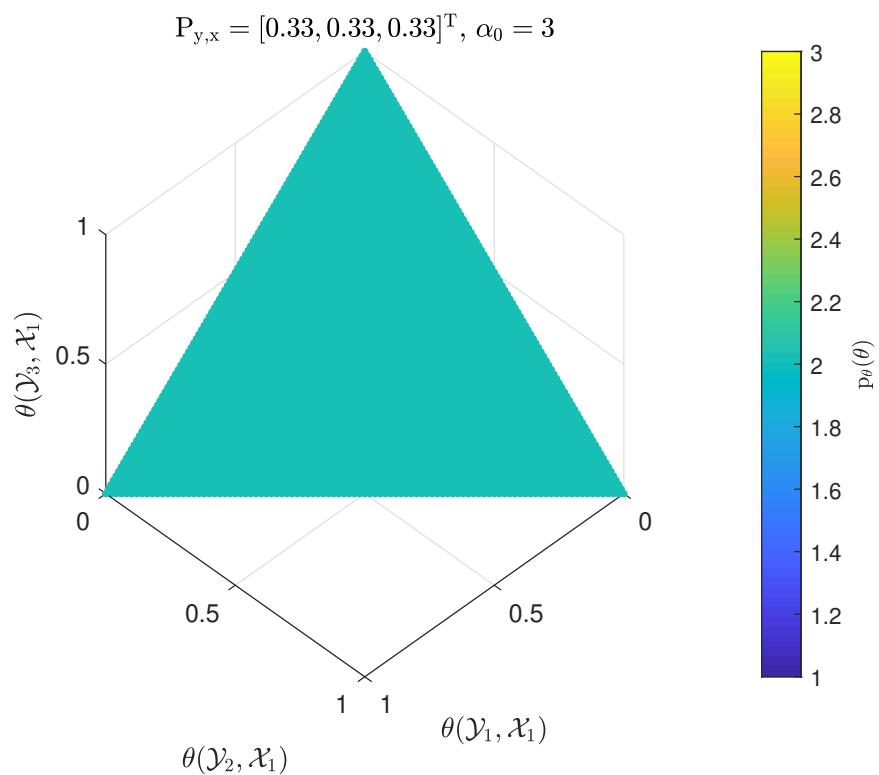


Figure 3.2: Uniform model prior PDF,  $|\mathcal{Y}| = 3, |\mathcal{X}| = 1$

PGR: implications of independence for posterior learning?

### 3.1.2 Training Set PMF, $P_D$

PGR: EVIDENCE TERMINOLOGY??

Next, the conditional distribution  $P_{D|\theta}$  will be used to determine the marginal PMF,  $P_D$  and properties will be discussed.

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 [11]

$$\begin{aligned} P_D(D) &= E_\theta \left[ \prod_{n=1}^N P_{D_n|\theta}(D_n|\theta) \right] \\ &= E_\theta \left[ \prod_{y \in \mathcal{Y}} \prod_{x \in \mathcal{X}} \theta(y, x)^{\bar{N}(y, x; D)} \right] \\ &= \beta(\alpha)^{-1} \beta(\alpha + \bar{N}(D)) . \end{aligned} \quad (3.10)$$

Note that values of the PMF  $P_D$  are equivalent to joint moments of the model  $\theta$ .

It is instructive to consider the limiting forms of this distribution for the extreme values of the model concentration parameter  $\alpha_0$ . As  $\alpha_0 \rightarrow \infty$ , the model concentrates at its mean and the training data  $D$  distribution is

$$\begin{aligned} P_D(D) &\rightarrow E_\theta \left[ \prod_{n=1}^N \theta(Y_n, X_n) \right] \\ &= \prod_{n=1}^N \frac{\alpha(Y_n, X_n)}{\alpha_0} . \end{aligned} \quad (3.11)$$

Conversely, as  $\alpha_0 \rightarrow 0$ , the distribution becomes

$$P_D(D) \rightarrow \sum_{y \in \mathcal{Y}} \sum_{x \in \mathcal{X}} \frac{\alpha(y, x)}{\alpha_0} \prod_{n=1}^N \delta[D_n, (y, x)] \quad (3.12)$$

and the training data are identical.

Next, the distribution of the sufficient statistic  $\bar{n}$  will be represented. As a Dirichlet distribution characterizes the parameters of the multinomial distribution  $P_{D|\theta}$ , the marginal PMF of  $\bar{n}$  is a Dirichlet-Multinomial distribution [9] parameterized by  $\alpha$ ,

$$\begin{aligned} P_{\bar{n}}(\bar{n}) &= \mathcal{M}(\bar{n}) \beta(\alpha)^{-1} \beta(\alpha + \bar{n}) \\ &= DM(\bar{n}; N, \alpha) . \end{aligned} \quad (3.13)$$

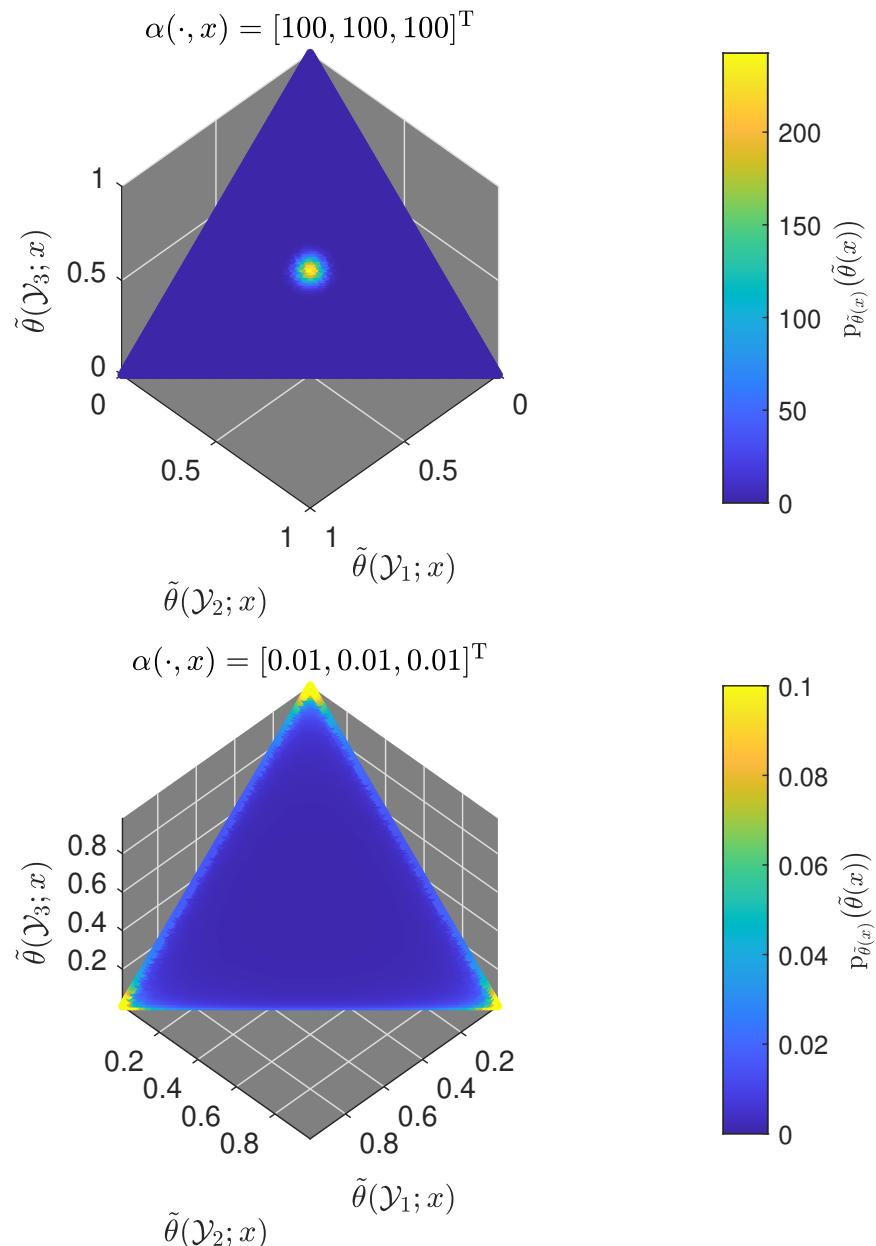


Figure 3.3: Model prior PDF for different concentrations  $\alpha'(x)$

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

$$\mu_{\bar{n}} = N \frac{\alpha}{\alpha_0} = N\mu_\theta \quad (3.14)$$

and

$$\begin{aligned} E_{\bar{n}} [\bar{n}(y, x)\bar{n}(y', x')] &= (3.15) \\ &= \frac{N}{\alpha_0(\alpha_0 + 1)} \left( (\alpha_0 + N)\alpha(y, x)\delta[y, y']\delta[x, x'] + (N - 1)\alpha(y, x)\alpha(y', x') \right) \\ &= \frac{N}{\alpha_0 + 1} \left( (\alpha_0 + N)\mu_\theta(y, x)\delta[y, y']\delta[x, x'] + \alpha_0(N - 1)\mu_\theta(y, x)\mu_\theta(y', x') \right). \end{aligned}$$

The covariance function is

$$\begin{aligned} \Sigma_{\bar{n}}(y, x, y', x') &= \frac{N(\alpha_0 + N)}{\alpha_0 + 1} (\mu_\theta(y, x)\delta[y, y']\delta[x, x'] - \mu_\theta(y, x)\mu_\theta(y', x')) \\ &= N(\alpha_0 + N)\Sigma_\theta(y, x, y', x'). \end{aligned} \quad (3.16)$$

Again, the data PMF's for minimal and maximal concentration  $\alpha_0$  are relevant. For  $\alpha_0 \rightarrow \infty$ , the model PDF  $p_\theta$  concentrates at its mean and thus  $\bar{n}$  is characterized by a multinomial distribution,

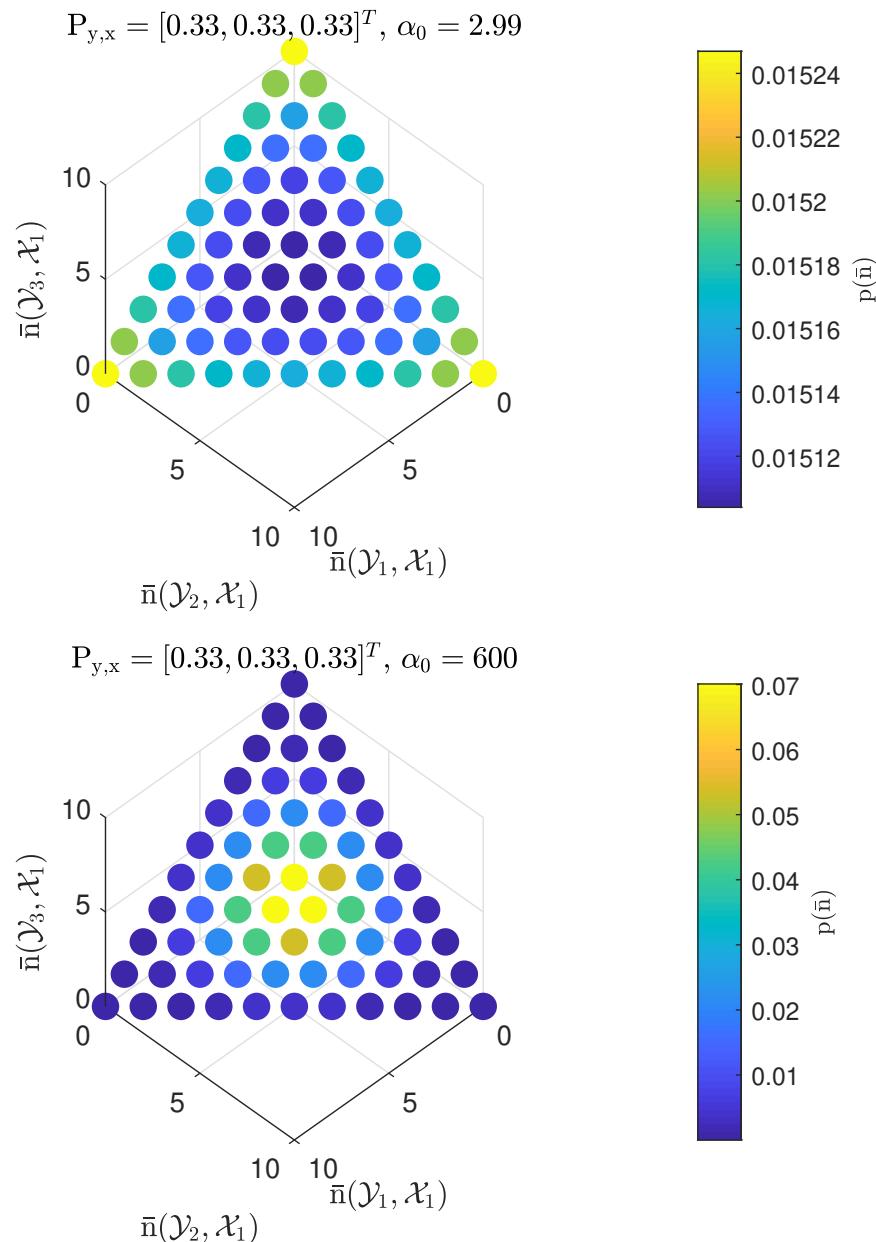
$$P_{\bar{n}}(\bar{n}) \rightarrow \mathcal{M}(\bar{n}) \prod_{y \in \mathcal{Y}} \prod_{x \in \mathcal{X}} \left( \frac{\alpha(y, x)}{\alpha_0} \right)^{\bar{n}(y, x)} \quad (3.17)$$

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

$$P_{\bar{n}}(\bar{n}) \rightarrow \sum_{y \in \mathcal{Y}} \sum_{x \in \mathcal{X}} \frac{\alpha(y, x)}{\alpha_0} \delta[\bar{n}, N\delta[\cdot, y]\delta[\cdot, x]]. \quad (3.18)$$

PGR: formal proofs for limiting PMFs? stirling/gautschi?

Figure 3.4 displays the distribution of  $\bar{n}$  for  $N = 10$  and different model concentrations  $\alpha_0$ . Observe that for large  $\alpha_0$ , the distribution approaches a multinomial distribution  $\bar{n} \sim \text{Multi}(N, \alpha/\alpha_0)$ . Figure 3.5 shows how a specific model prior influences the data PMF differently for different  $N$ . Observe that as the number of training samples increases, the PMF  $P_{\bar{n}}$  tends toward  $P_{\bar{n}}(\bar{n}) \approx N^{1-|\mathcal{Y}||\mathcal{X}|} p_\theta(\bar{n}/N)$ ; this can be proven using Gautschi's inequality [19].

Figure 3.4:  $P(\bar{n})$  for different prior concentrations  $\alpha_0$

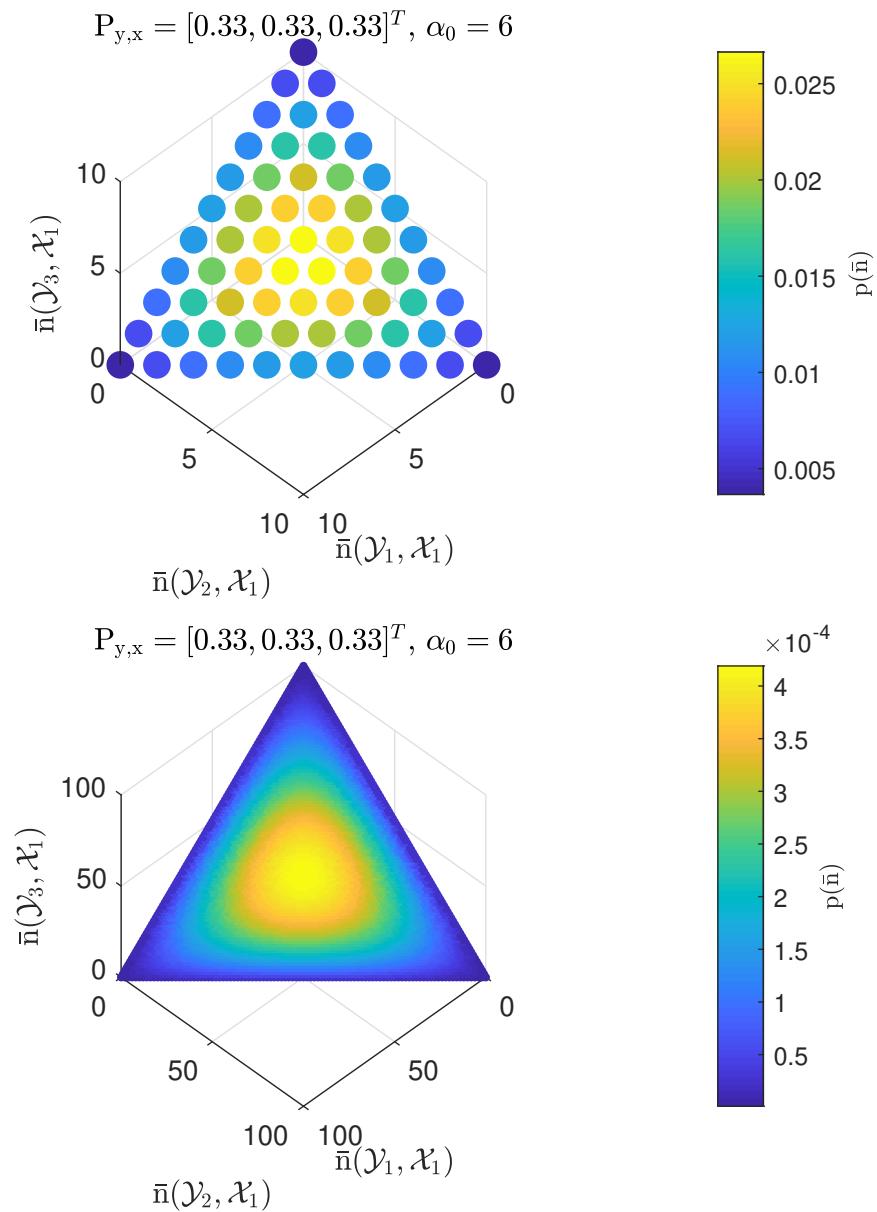
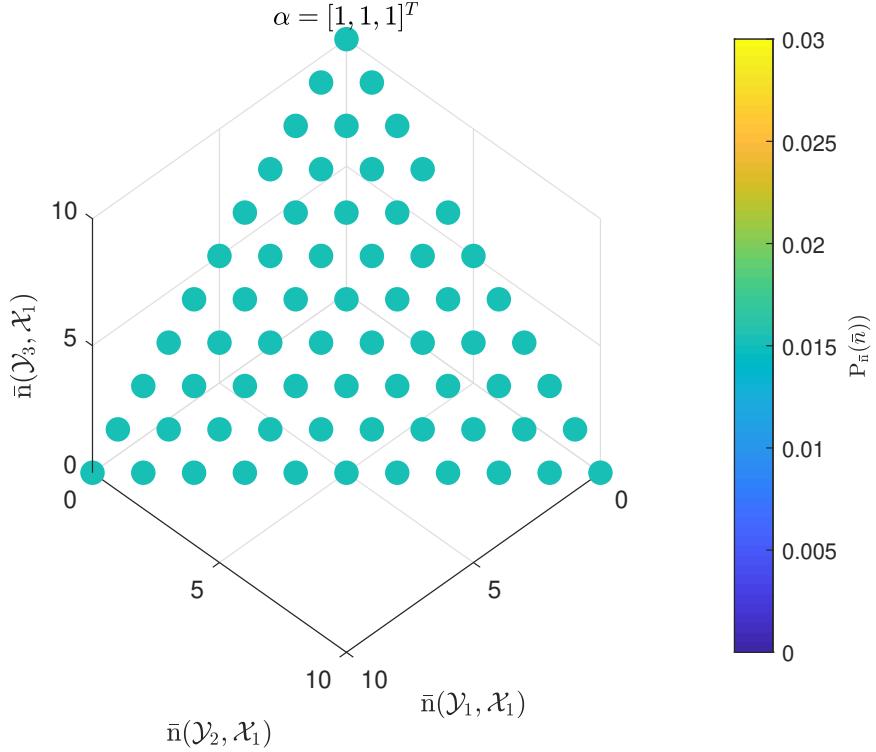


Figure 3.5:  $P(\bar{n})$  for different training set sizes  $N$

Figure 3.6:  $P(\bar{n})$  for uniform prior

**Uniform Prior** For the uniform distribution,  $\alpha(y, x) = 1$ ,

$$P_D(D) = \mathcal{M}(\{N, |\mathcal{Y}| |\mathcal{X}| - 1\})^{-1} \mathcal{M}(\bar{N}(D))^{-1} \quad (3.19)$$

and

$$P_{\bar{n}} = |\bar{\mathcal{N}}|^{-1} = \mathcal{M}(\{N, |\mathcal{Y}| |\mathcal{X}| - 1\})^{-1}. \quad (3.20)$$

The distribution of  $\bar{n}$  is uniform over the set  $\bar{\mathcal{N}}$ . The PMF for D depends on the training data only through the multinomial coefficient; consequently, more “concentrated” training sets are more probable.

### 3.1.2.1 Marginal and Conditional Distributions

It is also useful to express the marginal and conditional distributions for the training data given the Dirichlet prior. As  $P_{X|\theta}$  is of exponential form with respect to the marginal model

$\theta'$ , the marginal distribution of  $X$  can be expressed as

$$\begin{aligned} P_X(X) &= E_{\theta'} [P_{X|\theta'}](X) \\ &= E_{\theta} \left[ \prod_{n=1}^N P_{X_n|\theta}(X_n|\theta) \right] \\ &= E_{\theta'} \left[ \prod_{x \in \mathcal{X}} \theta'(x)^{N'(x;X)} \right] \\ &= \beta(\alpha')^{-1} \beta(\alpha' + N'(X)) . \end{aligned} \quad (3.21)$$

As the model marginal  $\theta'$  and conditional  $\tilde{\theta}$  are independent, the distribution  $P_{Y|X}$  can be represented as

$$\begin{aligned} P_{Y|X}(Y|X) &= E_{\tilde{\theta}} [P_{Y|X,\tilde{\theta}}](Y|X) \\ &= \prod_{x \in \mathcal{X}} E_{\tilde{\theta}(x)} \left[ \prod_{y \in \mathcal{Y}} \tilde{\theta}(y; x)^{\bar{N}(y,x;Y,X)} \right] \\ &= \prod_{x \in \mathcal{X}} \beta(\alpha(\cdot, x))^{-1} \beta(\alpha(\cdot, x) + \bar{N}(\cdot, x; Y, X)) . \end{aligned} \quad (3.22)$$

The corresponding distributions for the sufficient statistics will be expressed as well. Recall that  $n'|\theta \sim \text{Multi}(N, \theta')$ ; by the aggregation property of Dirichlet-Multinomial functions [9], the random process is distributed as  $n' \sim \text{DM}(N, \alpha')$ .

Also of interest is the distribution of  $\bar{n}$  conditioned on its aggregation  $n'$ . Using the Dirichlet-Multinomial properties presented in Appendix A.3, it can be shown that

$$\begin{aligned} P_{\bar{n}|n'}(\bar{n}|n') &= \prod_{x \in \mathcal{X}} \left[ \mathcal{M}(\bar{n}(\cdot, x)) \beta(\alpha(\cdot, x))^{-1} \beta(\alpha(\cdot, x) + \bar{n}(\cdot, x)) \right] \\ &= \prod_{x \in \mathcal{X}} \text{DM}(\bar{n}(\cdot, x); n'(x), \alpha(\cdot, x)) \end{aligned} \quad (3.23)$$

over the domain  $\left\{ \bar{n} \in \mathbb{Z}_{\geq 0}^{\mathcal{Y} \times \mathcal{X}} : \sum_{y \in \mathcal{Y}} \bar{n}(y, \cdot) = n' \right\}$ . 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, such that  $\bar{n}(\cdot, x)|n'(x) \sim \text{DM}(n'(x), \alpha(\cdot, x))$ .

### 3.1.3 Predictive PMF, $P_{y|x,D}$

As shown in Equation (2.8), the decision selected by the optimally designed function depends on  $P_{y|x,D}$ , the distribution of the unobserved  $y$  conditioned on all observable random elements. This PMF will be expressed next.

First observe that since  $P_{D|\theta}$  is of exponential form, the Dirichlet prior  $p_\theta$  is its conjugate prior [18]; thus, the model posterior PDF given the training data is

$$p_{\theta|D}(\theta|D) = \beta (\alpha + \bar{N}(D))^{-1} \prod_{y \in \mathcal{Y}} \prod_{x \in \mathcal{X}} \theta(y, x)^{\alpha(y, x) + \bar{N}(y, x; D) - 1}, \quad (3.24)$$

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

This posterior distribution is of specific interest in the machine learning literature. While Bayesian techniques are used here, often point estimates of the model  $\theta$  are formed; perhaps the most common approach is to form the Maximum a posteriori estimate,

$$\theta_{MAP}(D) = \arg \max_{\theta \in \Theta} P_{\theta|D}(\theta|D) = \frac{\bar{N}(D) + \alpha - 1}{N + \alpha_0 - |\mathcal{Y}| |\mathcal{X}|}. \quad (3.25)$$

This maximizing value is only valid when  $\bar{N}(D) > 1$ . For the uniform model prior, the maximizing value of the posterior is the empirical PMF  $\bar{N}(D)/N$ .

PGR: MAP discussion out of place??

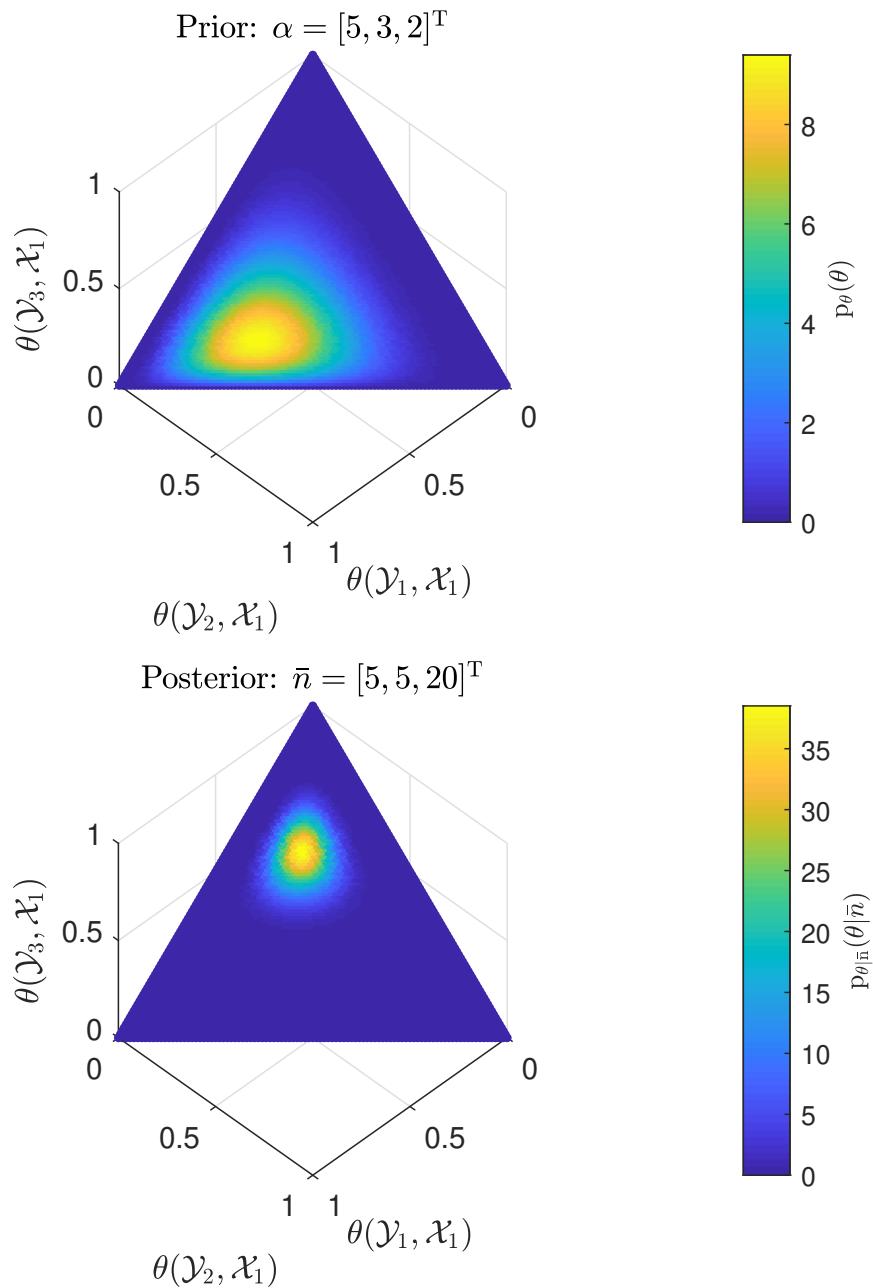
Also, the concentration parameter increases proportionately with increasing volumes of training data; consequently, as  $N \rightarrow \infty$ , the posterior converges to  $p_{\theta|D} \rightarrow \delta(\cdot - \bar{N}(D)/N)$ . Thus, as more data is collected, the model can be more positively identified and used to formulate minimum risk decisions. Conversely, as  $\alpha_0 \rightarrow \infty$ , the prior model certainty is stronger and the posterior tends toward  $p_{\theta|D} \rightarrow \delta(\cdot - \alpha/\alpha_0)$ , independent of the training data.

Figure 3.7 shows the influence of the training data on the model distribution; after conditioning on the training data (via  $\bar{n}$ ), the PDF concentration shifts away from the models favored by the prior knowledge and towards other models that better account for the observations.

The joint PMF of  $y$  and  $x$  conditioned on the training data is expressed as [12]

$$\begin{aligned} P_{y,x|D} &= \mu_{\theta|D} = \frac{\alpha + \bar{N}(D)}{\alpha_0 + N} \\ &= \left( \frac{\alpha_0}{\alpha_0 + N} \right) \frac{\alpha}{\alpha_0} + \left( \frac{N}{\alpha_0 + N} \right) \frac{\bar{N}(D)}{N} \\ &= \left( \frac{\alpha_0}{\alpha_0 + N} \right) P_{y,x} + \left( \frac{N}{\alpha_0 + N} \right) \frac{\bar{N}(D)}{N}. \end{aligned} \quad (3.26)$$

This is a mixture distribution of the prior expectation  $\mu_\theta = \alpha/\alpha_0$  and the empirical distribution  $\bar{N}(D)/N$ . The more informative the model prior (i.e. larger  $\alpha_0$ ), the more the

Figure 3.7: Model  $\theta$  PDF, prior and posterior

prior mean is favored; the more data, the more the empirical PMF is favored. The marginal distribution for  $x$  given  $D$  is

$$\begin{aligned} P_{x|D} &= \frac{\alpha' + N'(D)}{\alpha_0 + N} \\ &= \left( \frac{\alpha_0}{\alpha_0 + N} \right) \frac{\alpha'}{\alpha_0} + \left( \frac{N}{\alpha_0 + N} \right) \frac{N'(D)}{N} \\ &= \left( \frac{\alpha_0}{\alpha_0 + N} \right) P_x + \left( \frac{N}{\alpha_0 + N} \right) \frac{N'(D)}{N}. \end{aligned} \quad (3.27)$$

Finally, the predictive distribution of interest is generated via Bayes rule as

$$\begin{aligned} P_{y|x,D} &= \frac{\alpha(\cdot, x) + \bar{N}(\cdot, x; D)}{\alpha'(x) + N'(x; D)} \\ &= \left( \frac{\alpha'(x)}{\alpha'(x) + N'(x; D)} \right) \frac{\alpha(\cdot, x)}{\alpha'(x)} + \left( \frac{N'(x; D)}{\alpha'(x) + N'(x; D)} \right) \frac{\bar{N}(\cdot, x; D)}{N'(x; D)} \\ &= \left( \frac{\alpha'(x)}{\alpha'(x) + N'(x; D)} \right) P_{y|x} + \left( \frac{N'(x; D)}{\alpha'(x) + N'(x; D)} \right) \frac{\bar{N}(\cdot, x; D)}{N'(x; D)}. \end{aligned} \quad (3.28)$$

The last representation views the distribution as a convex combination of two conditional distributions. The first distribution  $P_{y|x} = \alpha(\cdot, x)/\alpha'(x)$  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 on  $D$ , not on  $\alpha$ .

The weighting factors  $\alpha'(x)$  and  $N'(x; D)$  are the concentration of the conditional prior  $\tilde{\theta}(x)$  and the number of training samples satisfying  $X_n = x$ . As  $N'(x; D)/\alpha'(x) \rightarrow 0$ , the PMF tends toward the conditional distribution  $P_{y|x}$ , which only depends on the model parameter  $\alpha$ . As  $N'(x; D)/\alpha'(x) \rightarrow \infty$ ,  $P_{y|x,D}$  tends towards the empirical conditional distribution.

**Uniform Prior** For the uniform model prior PDF, the conditional distribution is

$$\begin{aligned} P_{y|x,D} &= \frac{\bar{N}(\cdot, x; D) + 1}{N'(x; D) + |\mathcal{Y}|} \\ &= \left( \frac{|\mathcal{Y}|}{N'(x; D) + |\mathcal{Y}|} \right) \frac{1}{|\mathcal{Y}|} + \\ &\quad \left( \frac{N'(x; D)}{N'(x; D) + |\mathcal{Y}|} \right) \frac{\bar{N}(\cdot, x; D)}{N'(x; D)}. \end{aligned} \quad (3.29)$$

Now the prior PMF contribution  $P_{y|x}$  is a uniform distribution over the  $|\mathcal{Y}|$  possible outputs. The weighting factors are controlled by  $\alpha'(x) = |\mathcal{Y}|$ ; the more possible outcomes  $|\mathcal{Y}|$  there are for a given training set size, the more the conditional distribution tends toward the uniform PMF implied by the model prior.

### 3.1.3.1 Representation using the complete model posterior

PGR: reference posterior equations!

PGR: DIR FIGS? for PDF asymptotics?

The Bayesian distributions  $P_{x|D}$  and  $P_{y|x,D}$  can also be found from the posterior distributions  $p_{\theta'|D}$  and  $p_{\tilde{\theta}|x,D}$ , respectively. As the Dirichlet assumption renders  $\theta'$  and  $\tilde{\theta}$  independent, it can be shown that  $P_{Y|X} = E_{\tilde{\theta}} [P_{Y|X,\tilde{\theta}}]$  and thus that  $\theta'$  is conditionally independent of  $Y$  given  $X$ . Furthermore, the Dirichlet distribution  $p_{\theta'}$  is the conjugate prior for  $P_{X|\theta'}$ . As a result,  $\theta'|D \sim \text{Dir}(\alpha' + N'(X))$  and

$$\begin{aligned} P_{x|D}(D) &= \mu_{\theta'|D}(D) = \mu_{\theta'|X}(X) \\ &= \frac{\alpha' + N'(X)}{\alpha_0 + N}. \end{aligned} \quad (3.30)$$

Similarly, the distribution can be expressed in terms of the empirical PMF sufficient statistic as

$$\begin{aligned} P_{x|\bar{n}}(\bar{n}) &= \mu_{\theta'|\bar{n}}(\bar{n}) = \mu_{\theta'|n'} \left( \sum_y \bar{n}(y, \cdot) \right) \\ &= \frac{\alpha' + \sum_y \bar{n}(y, \cdot)}{\alpha_0 + N}, \end{aligned} \quad (3.31)$$

where the dependency on  $\bar{n}$  is expressed only through the marginal random process  $n'$ .

The posterior  $p_{\tilde{\theta}|x,D}$  can be simplified by noting that the independence of  $\theta'$  and  $\tilde{\theta}$  implies  $P_{Y|X,x} = E_{\tilde{\theta}} [P_{Y|X,\tilde{\theta}}] = P_{Y|X}$ . Consequently,  $\tilde{\theta}$  is conditionally independent of  $x$  given  $D$ . Thus, as  $p_{\tilde{\theta}}$  is a conjugate prior for  $P_{Y|X,\tilde{\theta}}$  the posterior distribution is

$$\begin{aligned} p_{\tilde{\theta}|D,x}(\tilde{\theta}|D, x) &= p_{\tilde{\theta}|D}(\tilde{\theta}|D) = \prod_{x' \in \mathcal{X}} p_{\tilde{\theta}(x')|D}(\tilde{\theta}(x')|D) \\ &= \prod_{x' \in \mathcal{X}} \text{Dir}(\tilde{\theta}(x'); \alpha(\cdot, x') + \bar{N}(\cdot, x'; D)) \end{aligned} \quad (3.32)$$

and the distinct model conditional PMF's are independent from one another. A similar treatment demonstrates that

$$\begin{aligned} p_{\tilde{\theta}|\bar{n},x}(\tilde{\theta}|\bar{n}, x) &= p_{\tilde{\theta}|\bar{n}}(\tilde{\theta}|\bar{n}) = \prod_{x' \in \mathcal{X}} p_{\tilde{\theta}(x')|\bar{n}(\cdot, x')}(\tilde{\theta}(x')|\bar{n}(\cdot, x')) \\ &= \prod_{x' \in \mathcal{X}} \text{Dir}(\tilde{\theta}(x'); \alpha(\cdot, x') + \bar{n}(\cdot, x')). \end{aligned} \quad (3.33)$$

Observe that when the conditioning is performed using the sufficient statistic, the independent conditional models  $\tilde{\theta}(x)$  are only dependent on their corresponding subset of the empirical PMF,  $\bar{n}(\cdot, x)$ .

The Bayes predictive PMF can thus be expressed as

$$\begin{aligned} P_{y|x,D}(x, D) &= \mu_{\tilde{\theta}(x)|x,D}(x, D) = \mu_{\tilde{\theta}(x)|D}(D) \\ &= \frac{\alpha(\cdot, x) + \bar{N}(\cdot, x; D)}{\alpha'(x) + N'(x; D)} \end{aligned} \quad (3.34)$$

or, via the sufficient statistic,

$$\begin{aligned} P_{y|x,\bar{n}}(x, \bar{n}) &= \mu_{\tilde{\theta}(x)|x,\bar{n}}(x, \bar{n}) = \mu_{\tilde{\theta}(x)|\bar{n}(\cdot,x)}(\bar{n}(\cdot, x)) \\ &= \frac{\alpha(\cdot, x) + \bar{n}(\cdot, x)}{\alpha'(x) + \sum_y \bar{n}(y, x)}. \end{aligned} \quad (3.35)$$

A consequence of the Dirichlet prior is that the predictive PMF for a given value of  $x$  only depends on the corresponding training data  $\bar{n}(\cdot, x)$ , such that  $P_{y|x,\bar{n}}(x, \bar{n}) = P_{y|x,\bar{n}(\cdot,x)}(x, \bar{n}(\cdot, x))$ . This is intuitive considering the independence of the conditional models  $\tilde{\theta}(x)$  from one another.

## 3.2 Model Estimation Perspective

PGR: generalize, move before Dirichlet?

PGR: use poster table??

It is instructive to treat the distribution  $P_{y|x,\bar{n}}$  as an estimate of the unknown conditional PMF  $P_{y|x,\theta} \equiv \tilde{\theta}(x)$  and investigate the effects of informative prior knowledge. For a given  $x$  and corresponding number of training samples  $n'(x)$ , the expected value of the estimate conditioned on the true model  $\theta$  is

$$\begin{aligned} E_{\bar{n}|n',\theta} [P_{y|x,\bar{n}}] &= E_{\bar{n}(\cdot,x)|n'(x),\tilde{\theta}(x)} [P_{y|x,\bar{n}(\cdot,x)}] \\ &= \left( \frac{\alpha'(x)}{\alpha'(x) + n'(x)} \right) \frac{\alpha(\cdot, x)}{\alpha'(x)} + \left( \frac{n'(x)}{\alpha'(x) + n'(x)} \right) \tilde{\theta}(x), \end{aligned} \quad (3.36)$$

where the properties of a multinomial distribution conditioned on its aggregation have been used. The result is a convex combination of the conditional data-independent distribution

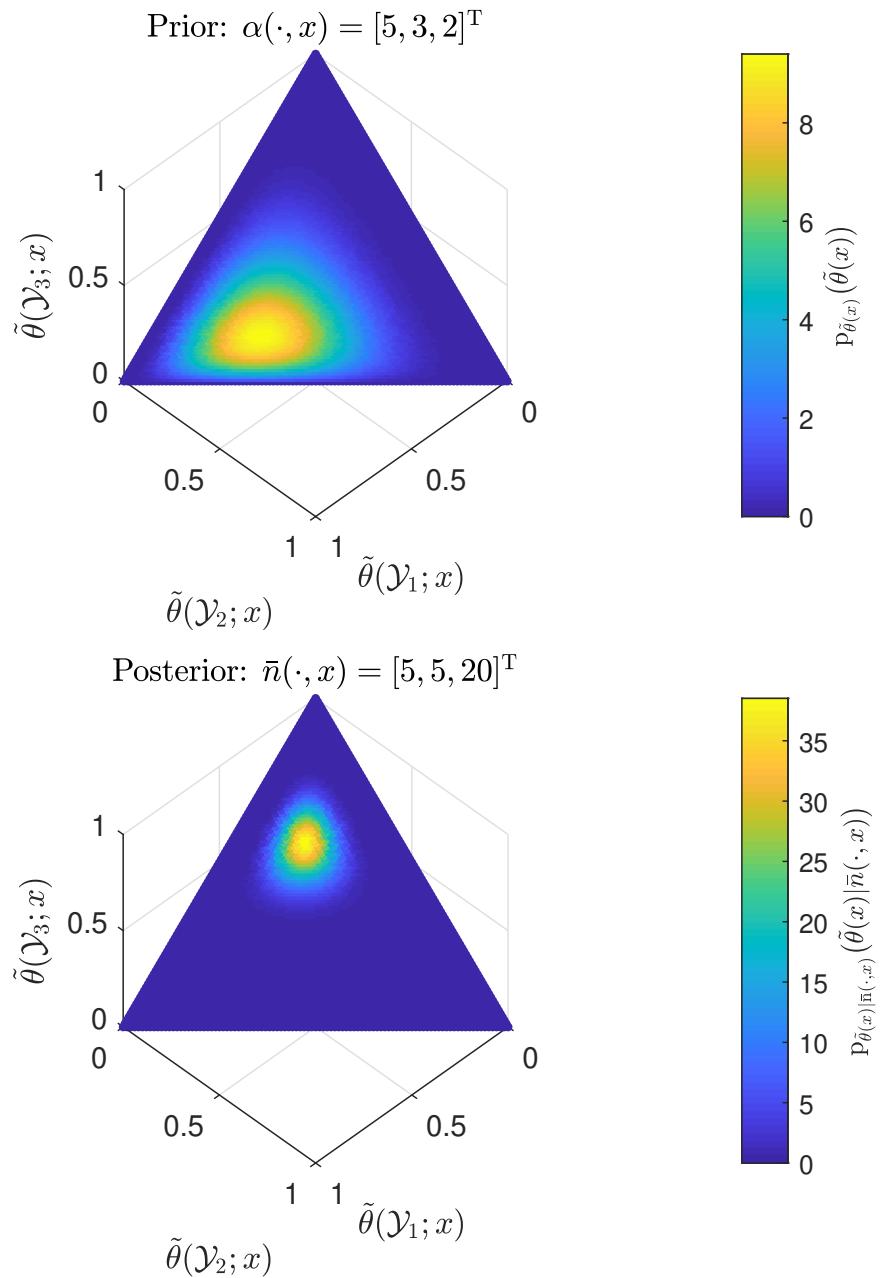


Figure 3.8: Model PDF, prior and posterior

$\alpha(\cdot, x)/\alpha'(x)$  and the true conditional distribution  $\tilde{\theta}(x)$ . The convex coefficients are dependent on the “marginal” values  $\alpha'$  and  $n'$ ; note that as the number of matching training samples  $n'(x)$  increases relative to  $\alpha'(x)$ , the estimate tends towards the true conditional PMF.

PGR: suppress delta dependency on nbar and theta?

To aid characterization of the estimator, define the random process  $\Delta(x, \bar{n}, \theta) \equiv P_{y|x,\bar{n}} - P_{y|x,\theta} \in \mathbb{R}^{\mathcal{Y}}$ . For a given  $x$  and corresponding number of training samples  $n'(x)$ , the bias of the conditional PMF estimate is

$$\begin{aligned}\text{Bias}(x, n', \theta) &= E_{\bar{n}|n',\theta} [\Delta(x, \bar{n}, \theta)] \\ &= \frac{\alpha'(x)}{\alpha'(x) + n'(x)} \left( \frac{\alpha(\cdot, x)}{\alpha'(x)} - \tilde{\theta}(x) \right)\end{aligned}\tag{3.37}$$

and its covariance function is

$$\begin{aligned}\text{Cov}(y, y'; x, n', \theta) &= C_{\bar{n}|n',\theta} [P_{y|x,\bar{n}}(\cdot|x, \bar{n})] (y, y') \\ &= \frac{\sum_{\bar{n}(\cdot,x)|n'(x),\tilde{\theta}(x)} (y, y')}{(\alpha'(x) + n'(x))^2} \\ &= \frac{n'(x)}{(\alpha'(x) + n'(x))^2} \left( \tilde{\theta}(y; x) \delta[y, y'] - \tilde{\theta}(y; x) \tilde{\theta}(y'; x) \right),\end{aligned}\tag{3.38}$$

where the properties of multinomial random processes have been used. Note that the bias is proportionate to the difference between the true conditional model and the data-independent estimate. The scaling factor tends from one to zero as  $n'(x)/\alpha'(x)$  tends from zero to infinity; as such, more informative priors (large  $\alpha'(x)$ ) will lead to PMF estimates that are prone to bias. Conversely, the variance of the PMF estimate tends to zero as  $\alpha'(x) \rightarrow \infty$ .

Combining the estimator bias and variance, the conditional second moments of  $\Delta(x, \bar{n}, \theta)$  are

$$\begin{aligned}\mathcal{E}(y, y'; x, n', \theta) &= E_{\bar{n}|n',\theta} [\Delta(y; x, \bar{n}, \theta) \Delta(y'; x, \bar{n}, \theta)] \\ &= \text{Bias}(y; x, n', \theta) \text{Bias}(y'; x, n', \theta) + \text{Cov}(y, y'; x, n', \theta).\end{aligned}\tag{3.39}$$

As  $n'(x) \rightarrow \infty$ , this function tends to zero and thus the underlying model  $\tilde{\theta}(x)$  is determined precisely. A more practical case is estimation with a finite volume of training data. Specification of the Dirichlet model prior can be interpreted as providing a distribution estimate  $\alpha(\cdot, x)/\alpha'(x)$  and a confidence level  $\alpha'(x)$ . Higher confidence reduces error due to

the variance of the estimator, but increases the error due to bias between the true model and its estimate; low confidence renders the estimate unbiased, but maximizes the estimator variance.

Also of interest, the conditional expectation of  $\mathcal{E}(\cdot, \cdot; x, n', \theta)$  is

$$\begin{aligned} & E_{x,n'|\theta} [\mathcal{E}(y, y'; x, n', \theta)] \\ &= E_{x|\theta} \left[ E_{n'(x)|\theta} \left[ \left( \frac{\alpha'(x)}{\alpha'(x) + n'(x)} \right)^2 \right] \left( \frac{\alpha(y, x)}{\alpha'(x)} - \tilde{\theta}(y; x) \right) \left( \frac{\alpha(y', x)}{\alpha'(x)} - \tilde{\theta}(y'; x) \right) \right] \\ &+ E_{x|\theta} \left[ E_{n'(x)|\theta} \left[ \frac{n'(x)}{(\alpha'(x) + n'(x))^2} \right] \left( \tilde{\theta}(y; x) \delta[y, y'] - \tilde{\theta}(y; x) \tilde{\theta}(y'; x) \right) \right]. \end{aligned} \quad (3.40)$$

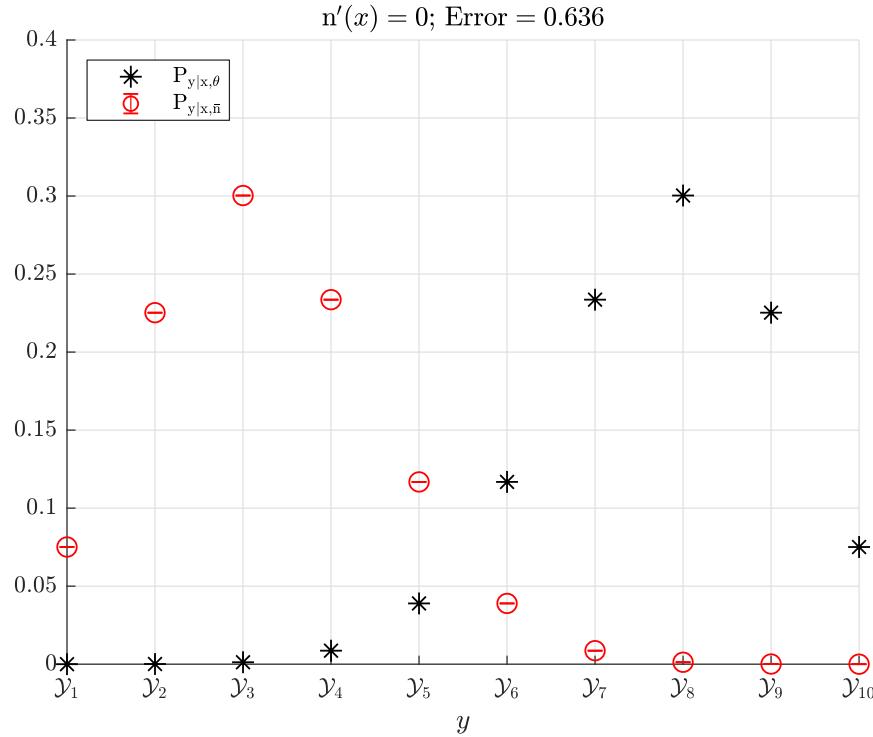
To exemplify how the model estimate  $P_{y|x,\bar{n}}$  approximates  $P_{y|x,\theta}$ , consider a scenario with  $|\mathcal{Y}| = 10$ . The data-independent PMF  $\alpha(\cdot, x)/\alpha'(x)$  and true model  $\tilde{\theta}(x)$  are shown in Figure 3.9 - note the significant mismatch.

PGR: ABOVE - STATE XCAL cardinality = 1

Figures 3.10 and 3.11 show how the bias and variance of the estimate change for different values of  $n'(x)$  and  $\alpha'(x)$ . The plot markers represent the conditional mean of the estimator,  $E_{\bar{n}|n',\theta} [P_{y|x,\bar{n}}(y|x, \bar{n})]$ ; the upper and lower error bars represent the square-root of the expected squared deviation above and below the conditional mean, respectively. Each individual plot heading provides the error  $\sqrt{\sum_{y \in \mathcal{Y}} \mathcal{E}(y, y; x, n', \theta)}$  to assess the quality of the PMF estimate.

Observe that for  $n'(x) = 1$ , the high variance of the  $\alpha'(x) = 0.1$  estimate (favoring the empirical PMF) renders it worse than the  $\alpha_0 = 10$  estimate; in fact, the variance is so high that the error exceeds that of the data-independent estimate  $\alpha(\cdot, x)/\alpha'(x)$  (Figure 3.9). Conversely, for  $n'(x) = 10$ , the confidence of the  $\alpha'(x) = 10$  estimate leads to high bias and the  $\alpha'(x) = 0.1$  estimate is superior. For  $n'(x) = 100$ , both the  $\alpha'(x) = 0.1$  and  $\alpha'(x) = 10$  estimates begin converging to the true distribution - this is guaranteed due to the full support of the Dirichlet prior.

PGR: full support discussion?

Figure 3.9: Model  $\theta$  estimate, no training data

### 3.3 Applications to Common Loss Functions

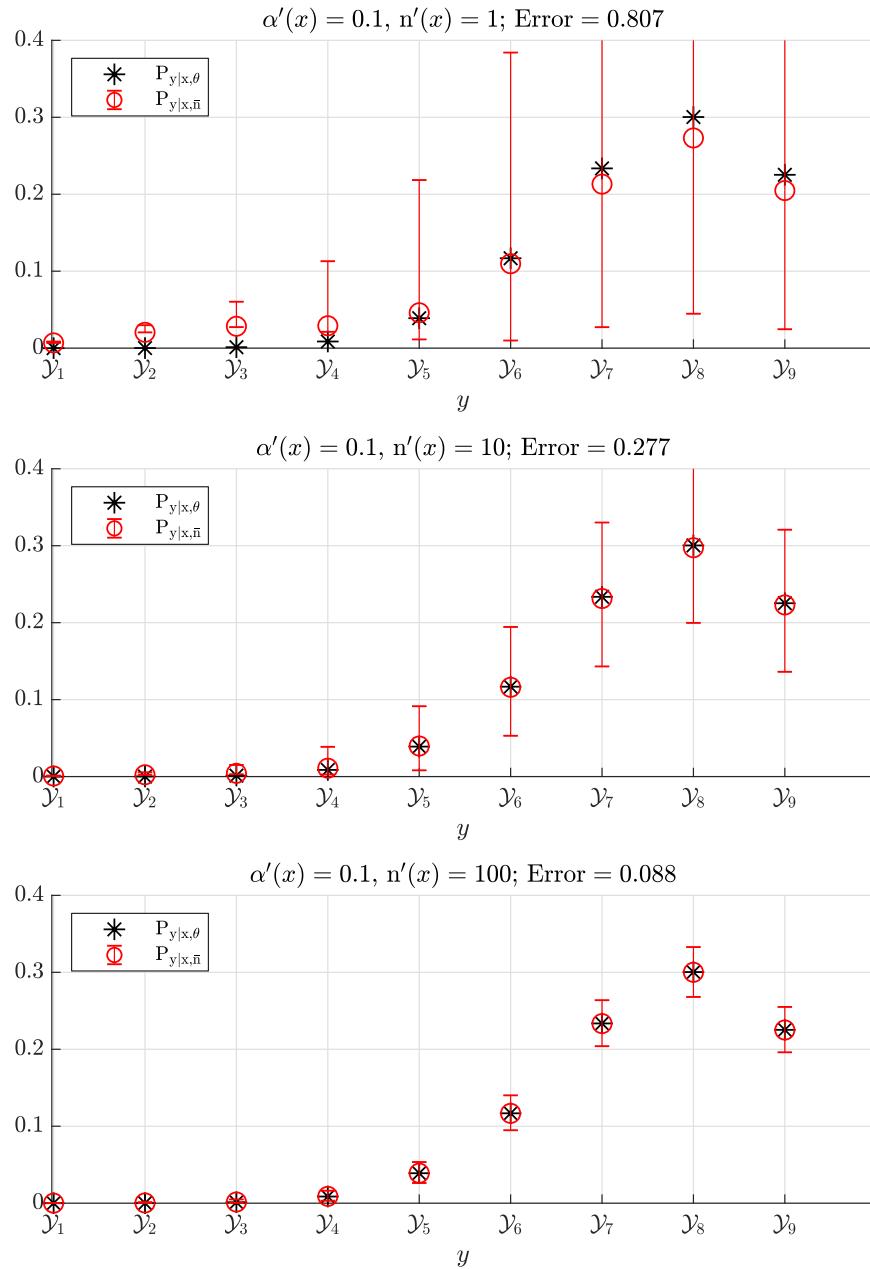
PGR: REMOVE GENERAL, RELOCATED MATERIAL!

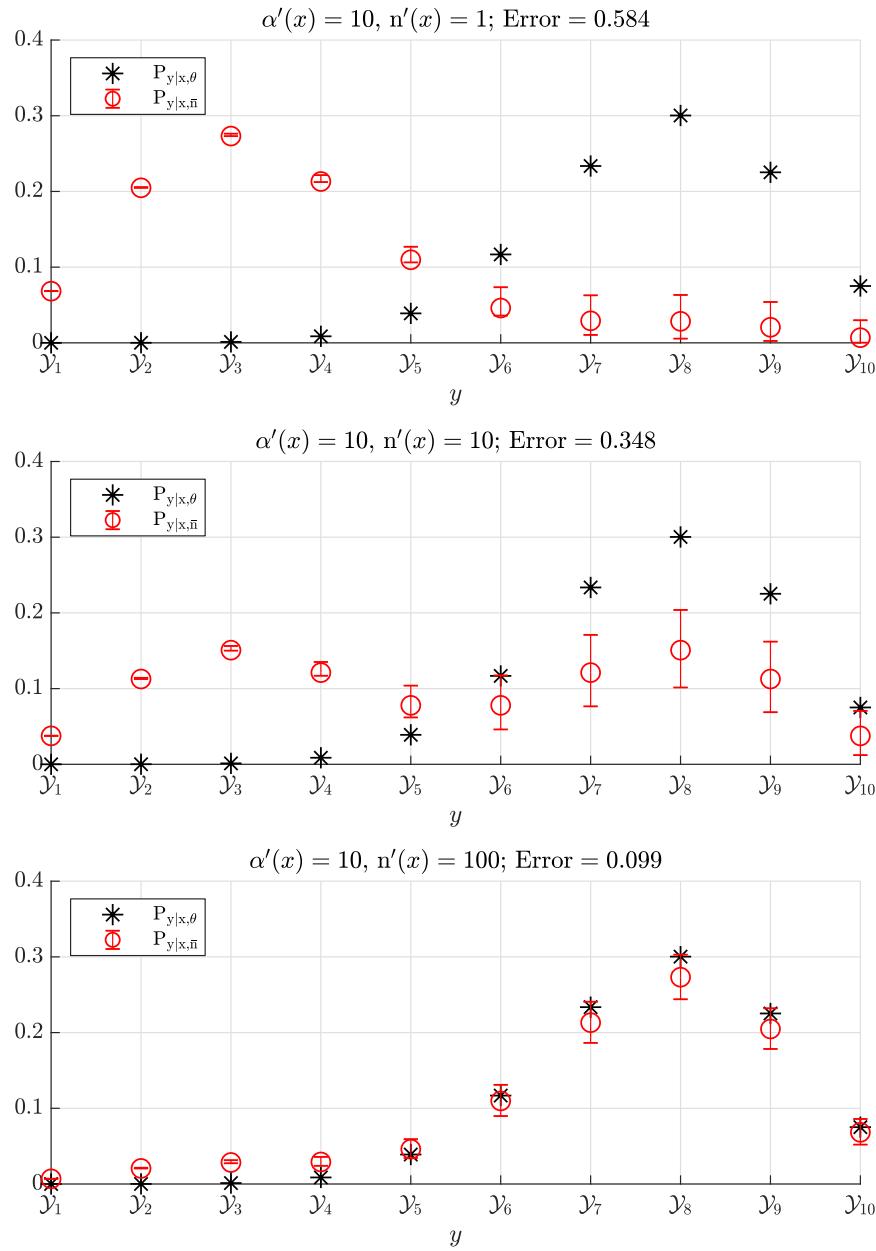
PGR: equations, plots for specific theta results? informative/non-informative tradeoff??? alpha/theta mismatch results???

In this section, the Dirichlet prior is applied to the regression and classification applications. Optimal learners  $f^*$  are found, the corresponding minimum Bayes risk  $\mathcal{R}^*$  is assessed, and the conditional risk  $\mathcal{R}_\Theta(f^*; \theta)$  is analyzed.

PGR: add formula for  $f(D)$ ? EMPIRICAL RISK DISCUSS, REGULARIZING weight

It is useful to substitute the Bayes predictive distribution using the Dirichlet prior (3.28)

Figure 3.10: Model  $\theta$  estimates,  $\alpha_0 = 0.1$

Figure 3.11: Model  $\theta$  estimates,  $\alpha_0 = 10$

into Equation (2.8), expressing the decision for a given input  $x$  and training set  $D$  as

$$\begin{aligned}
 f^*(x; D) &= \arg \min_{h \in \mathcal{H}} E_{y|x,D} [\mathcal{L}(h, y)] \\
 &= \arg \min_{h \in \mathcal{H}} \frac{\sum_{y \in \mathcal{Y}} \alpha(y, x) \mathcal{L}(h, y) + \sum_{y \in \mathcal{Y}} \bar{N}(y, x; D) \mathcal{L}(h, y)}{\alpha'(x) + N'(x; D)} \\
 &= \arg \min_{h \in \mathcal{H}} \left( \frac{\alpha'(x)}{\alpha'(x) + N'(x; D)} \right) \sum_{y \in \mathcal{Y}} \frac{\alpha(y, x)}{\alpha'(x)} \mathcal{L}(h, y) + \left( \frac{N'(x; D)}{\alpha'(x) + N'(x; D)} \right) \sum_{y \in \mathcal{Y}} \frac{\bar{N}(y, x; D)}{N'(x; D)} \mathcal{L}(h, y) \\
 &= \arg \min_{h \in \mathcal{H}} \left( \frac{\alpha'(x)}{\alpha'(x) + N'(x; D)} \right) E_{y|x} [\mathcal{L}(h, y)] + \left( \frac{N'(x; D)}{\alpha'(x) + N'(x; D)} \right) \frac{\sum_{n=1}^N \delta[x, X_n] \mathcal{L}(h, Y_n)}{\sum_{n=1}^N \delta[x, X_n]}.
 \end{aligned} \tag{3.41}$$

The metric to be minimized can be represented as a convex combination of two expected losses. The first expected loss is evaluated with respect to the conditional distribution  $P_{y|x} = \alpha(\cdot, x)/\alpha'(x)$ , which reflects the prior knowledge imparted by the model parameter  $\alpha$ . The second term is the conditional empirical risk, or the average loss among samples  $Y_n$  whose corresponding values  $X_n$  match the observed value  $x$ . The convex weights are inherited from the conditional distribution  $P_{y|x,D}$ ; thus, for a given observation  $x$ , the model prior parameter  $\alpha'(x)$  and the number of matching training samples  $N'(x; D)$  dictate which of the two expectations are emphasized.

### 3.3.1 Regression: the Squared-Error Loss

PGR: Use finite hypothesis space instead, wait for continuous DP???

PGR: add Dir conditional risk and analysis!!!

The elements of the finite cardinality set  $\mathcal{Y}$  are real numbers, such that  $\mathcal{Y} \subset \mathbb{R}$ . Again,  $\mathcal{H} = \mathbb{R} \supset \mathcal{Y}$ .

#### 3.3.1.1 Optimal Estimate: the Posterior Mean

PGR: plots?

Substituting in the Bayes predictive distribution for a Dirichlet prior (3.28) into (2.31),

the optimal Bayesian estimate is

$$\begin{aligned} f^*(x; D) &= \mu_{y|x,D} \\ &= \left( \frac{\alpha'(x)}{\alpha'(x) + N'(x; D)} \right) \sum_{y \in \mathcal{Y}} y \frac{\alpha(y, x)}{\alpha'(x)} + \left( \frac{N'(x; D)}{\alpha'(x) + N'(x; D)} \right) \sum_{y \in \mathcal{Y}} y \frac{\bar{N}(y, x; D)}{N'(x; D)} \\ &= \left( \frac{\alpha'(x)}{\alpha'(x) + N'(x; D)} \right) \mu_{y|x} + \left( \frac{N'(x; D)}{\alpha'(x) + N'(x; D)} \right) \frac{\sum_{n=1}^N \delta[x, X_n] Y_n}{N'(x; D)}. \end{aligned} \quad (3.42)$$

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 mean of the training values  $Y_n$  which have a value  $X_n$  matching the observed value  $x$ . 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.

**Uniform Prior** The optimal estimator for a uniform prior is

$$\begin{aligned} f^*(x; D) &= \left( \frac{|\mathcal{Y}|}{N'(x; D) + |\mathcal{Y}|} \right) \frac{1}{|\mathcal{Y}|} \sum_{y \in \mathcal{Y}} y + \left( \frac{N'(x; D)}{N'(x; D) + |\mathcal{Y}|} \right) \sum_{y \in \mathcal{Y}} y \frac{\bar{N}(y, x; D)}{N'(x; D)} \\ &= \left( \frac{|\mathcal{Y}|}{N'(x; D) + |\mathcal{Y}|} \right) \frac{1}{|\mathcal{Y}|} \sum_{y \in \mathcal{Y}} y + \left( \frac{N'(x; D)}{N'(x; D) + |\mathcal{Y}|} \right) \frac{\sum_{n=1}^N \delta[x, X_n] Y_n}{N'(x; D)}. \end{aligned} \quad (3.43)$$

Now, the model prior contribution to the weighting factors depends on the cardinality  $|\mathcal{Y}|$  and the prior expectation is simply the average of the elements of  $\mathcal{Y}$ .

### 3.3.1.2 Minimum Risk: the Expected Posterior Variance

PGR: determine irreducible risk separately, before??

The minimum Bayes squared-error is  $\mathcal{R}^* = E_{x,D} [\Sigma_{y|x,D}]$ . Using the sufficient statistic  $\bar{n} \equiv \bar{N}(D)$ , the minimum risk can also be represented as  $E_{x,\bar{n}} [\Sigma_{y|x,\bar{n}}]$ ; as such, the expectations are performed over  $\bar{n}$ . Decompose the conditional variance as

$$\Sigma_{y|x,\bar{n}} = E_{y|x,\bar{n}}[y^2] - \mu_{y|x,\bar{n}}^2 \quad (3.44)$$

and assess the expected values of these terms separately using distributions derived from the

Dirichlet prior. The first term is simply

$$\begin{aligned} E_{x,\bar{n}} [E_{y|x,\bar{n}}[y^2]] &= E_y[y^2] = \sum_{y \in \mathcal{Y}} y^2 \left( \sum_{x \in \mathcal{X}} \frac{\alpha(y, x)}{\alpha_0} \right) \\ &= E_x [E_{y|x}[y^2]] = \sum_{x \in \mathcal{X}} \frac{\alpha'(x)}{\alpha_0} \sum_{y \in \mathcal{Y}} y^2 \frac{\alpha(y, x)}{\alpha'(x)}, \end{aligned}$$

where the different functions of  $\alpha$  are represented by the PMF's of  $y$  and  $x$ . Next, find,

$$\begin{aligned} E_{x,\bar{n}} [\mu_{y|x,\bar{n}}^2] &= E_x \left[ E_{\bar{n}|x} \left[ \frac{(\alpha'(x)\mu_{y|x} + \sum_{y \in \mathcal{Y}} y\bar{n}(y,x))^2}{\alpha'(x)(\alpha'(x) + n'(x))^2} \right] \right] \\ &= E_x \left[ E_{\bar{n}} \left[ \frac{\alpha_0(\alpha'(x)\mu_{y|x} + \sum_{y \in \mathcal{Y}} y\bar{n}(y,x))^2}{\alpha'(x)(\alpha'(x) + n'(x))(\alpha_0 + N)} \right] \right] \\ &= E_x \left[ E_{n'} \left[ \frac{\alpha_0 E_{\bar{n}|n'} \left[ (\alpha'(x)\mu_{y|x} + \sum_{y \in \mathcal{Y}} y\bar{n}(y,x))^2 \right]}{\alpha'(x)(\alpha'(x) + n'(x))(\alpha_0 + N)} \right] \right] \\ &= \dots \\ &= E_x \left[ \frac{\alpha_0 E_{n'} \left[ n'(x) E_{y|x}[y^2] + (\alpha'(x) + n'(x) + 1)\alpha'(x)\mu_{y|x}^2 \right]}{\alpha'(x)(\alpha'(x) + 1)(\alpha_0 + N)} \right] \\ &= E_x \left[ \frac{N E_{y|x}[y^2] + (\alpha_0\alpha'(x) + N\alpha'(x) + \alpha_0)\mu_{y|x}^2}{(\alpha'(x) + 1)(\alpha_0 + N)} \right]. \end{aligned} \tag{3.45}$$

PGR: provide additional steps?

The above formulation exploits 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  $\bar{n}(\cdot, x) | n'(x) \sim \text{DM}(n'(x), \alpha(\cdot, x))$ .

PGR: move to appendix???

Finally, combine the two formulas to represent the minimum Bayes risk,

$$\begin{aligned} \mathcal{R}^* &= E_{x,\bar{n}} [E_{y|x,\bar{n}}[y^2] - \mu_{y|x,\bar{n}}^2] \\ &= E_x \left[ \frac{\alpha_0\alpha'(x) + N\alpha'(x) + \alpha_0}{(\alpha'(x) + 1)(\alpha_0 + N)} \Sigma_{y|x} \right] \\ &= E_x \left[ \frac{P_x(x) + (\alpha_0 + N)^{-1}}{P_x(x) + \alpha_0^{-1}} \Sigma_{y|x} \right]. \end{aligned} \tag{3.46}$$

The minimum risk is the expected value of the scaled conditional variance with respect to  $P_{y|x} = \alpha(\cdot, x)/\alpha'(x)$ . The expectation is taken with respect to the prior marginal distribution  $P_x = \alpha'/\alpha_0$ .

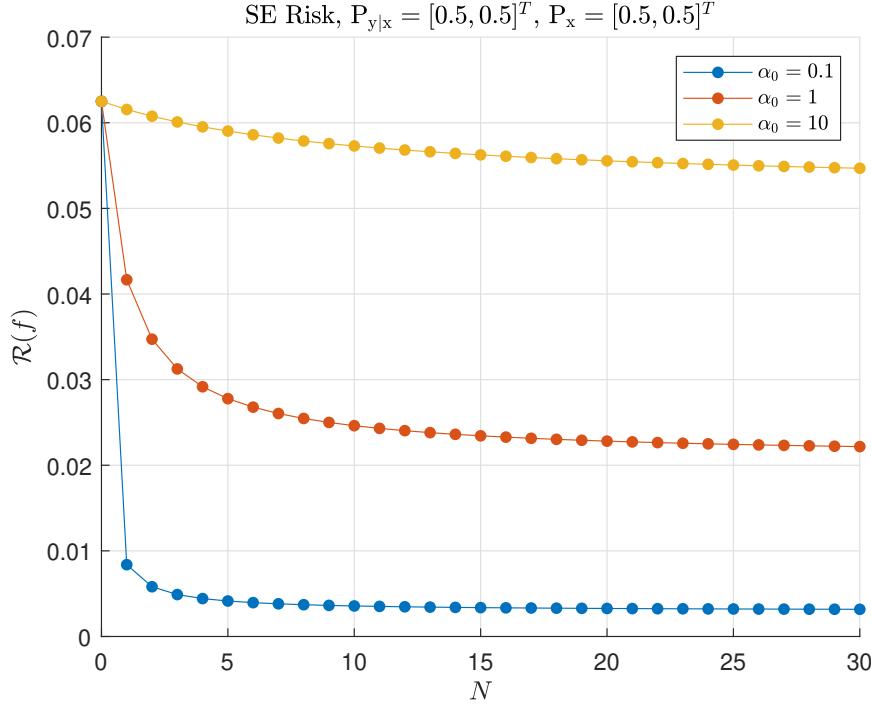


Figure 3.12: Minimum SE Risk for different training set sizes  $N$

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

PGR: first/second derivatives of alpha0??

To illustrate these trends, explicitly define the sets  $\mathcal{Y} = \{i/M_y : i = 0, \dots, M_y - 1\}$  and  $\mathcal{X} = \{i/M_x : i = 0, \dots, M_x - 1\}$ . Assume that the conditional variance  $\Sigma_{y|x}$  is independent of  $x$ ; in this case, the squared-error becomes the conditional variance scaled by a factor dependent on the marginal distribution  $P_x$ , such that  $\mathcal{R}^* = \Sigma_{y|x} E_x \left[ \frac{P_x(x) + (\alpha_0 + N)^{-1}}{P_x(x) + \alpha_0^{-1}} \right]$ . Figures 3.12 and 3.13 display how the risk changes with  $N$  and  $\alpha_0$  when  $P_{y|x}$  and  $P_x$  are fixed.

It may not seem intuitive for the risk to decrease when  $\alpha_0$  is smaller – the variance of the model  $\theta$  increases and the prior knowledge is less definitive. This is a result of the Dirichlet PDF weight shifting towards the  $|\mathcal{Y}| |\mathcal{X}|$  models which have  $\ell_0$  norms satisfying

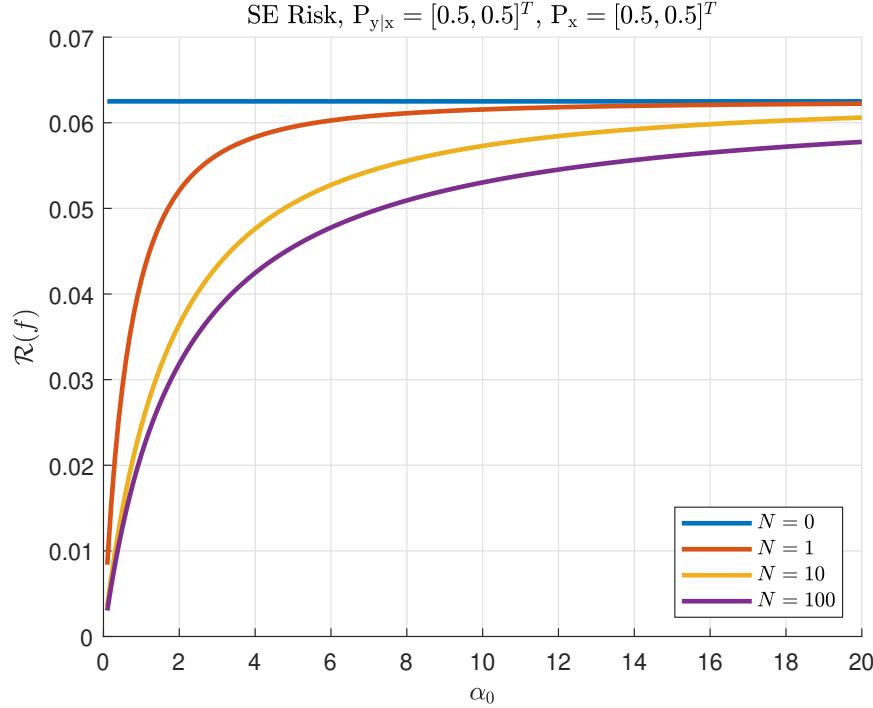
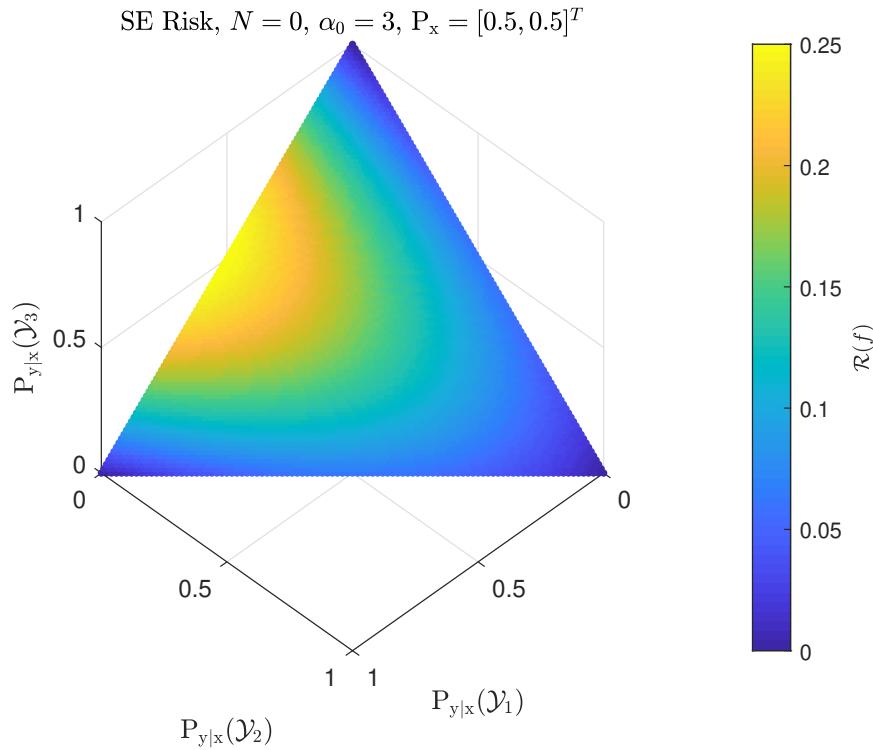
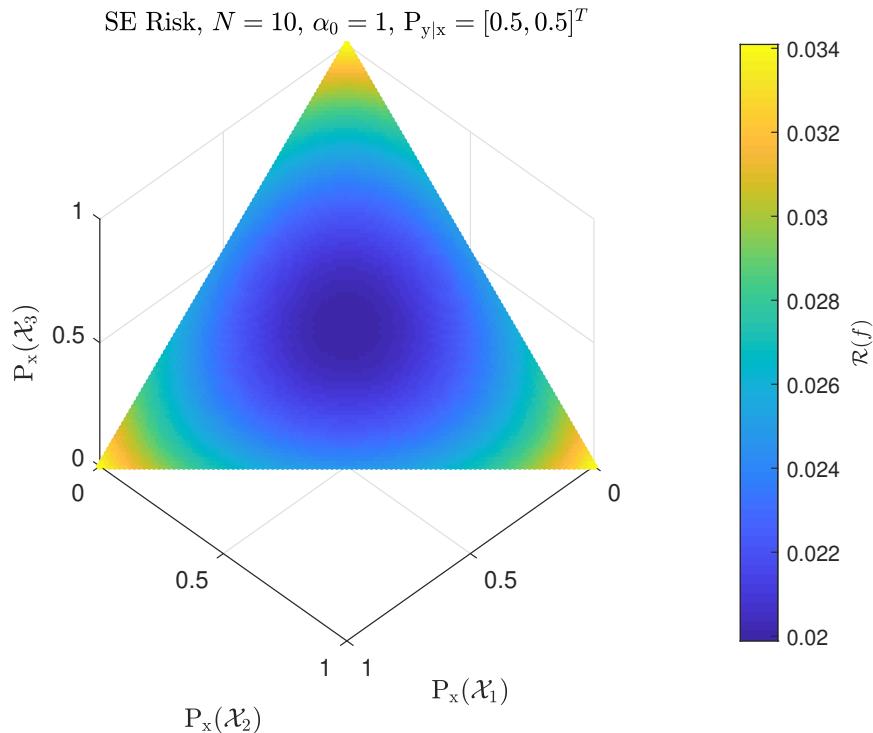


Figure 3.13: Minimum SE Risk for different prior concentrations  $\alpha_0$

$\|\theta\|_0 = 1$ . Although these PMF's are maximally separated (and uncorrelated), they all have zero variance. The optimal learner (3.42) will simply use the empirical distribution supplied via the training data - this allows exact identification of  $\theta$  with a single training pair.

It is also instructional to visualize how the minimum squared-error changes for fixed volume of training data  $N$  and a fixed prior concentration  $\alpha_0$ . First, consider how the risk changes with the conditional PMF  $P_{y|x}$ . Figure 3.14 demonstrates how the squared-error tends towards zero for PMFs that have  $\ell_0$ -norm equal to one.

Next, consider the effect of the marginal distribution  $P_x$ . Figure 3.15 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. Figures 3.16 and 3.17 show how different marginals  $P_x$  affect the risk as a function of  $N$  and  $\alpha_0$ , respectively.

Figure 3.14: Minimum SE Risk for different prior means  $P_{y|x}$ Figure 3.15: Minimum SE Risk for different prior means  $P_x$

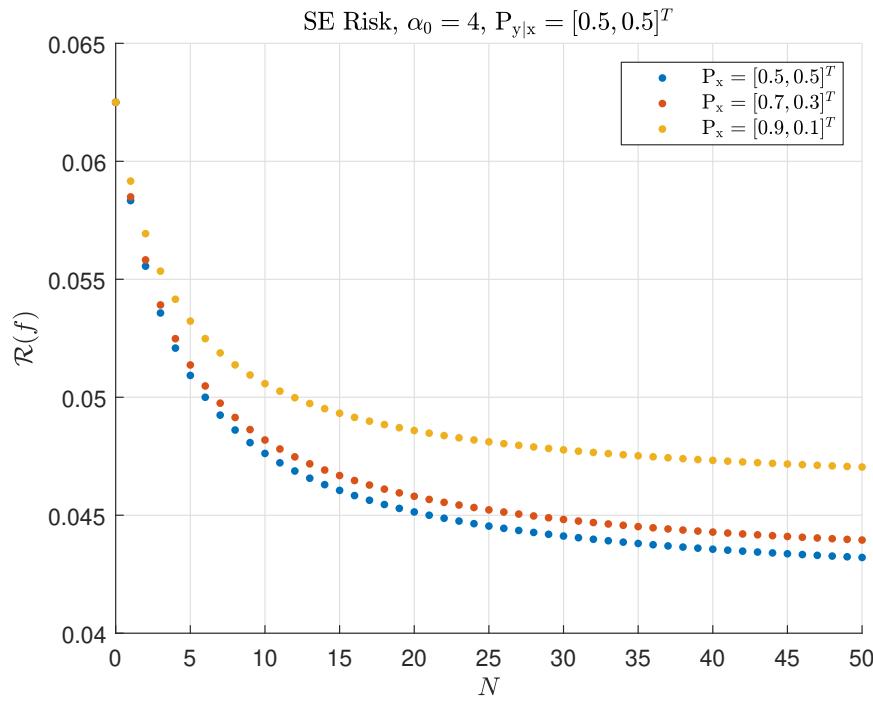


Figure 3.16: Minimum SE Risk for different training set volumes  $N$

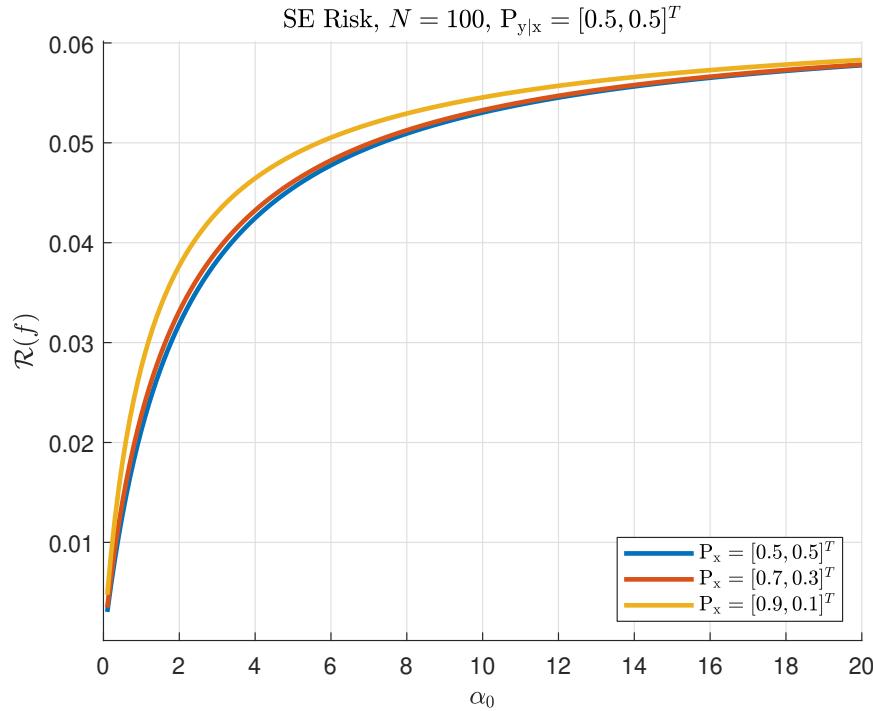


Figure 3.17: Minimum SE Risk for different prior concentrations  $\alpha_0$

**Uniform Prior** For the uniform model prior, the risk reduces to

$$\begin{aligned}\mathcal{R}^* &= \frac{|\mathcal{Y}|(N/|\mathcal{X}| + |\mathcal{Y}| + 1)}{(|\mathcal{Y}| + 1)(N/|\mathcal{X}| + |\mathcal{Y}|)} \left[ \left( \frac{1}{|\mathcal{Y}|} \sum_{y \in \mathcal{Y}} y^2 \right) - \left( \frac{1}{|\mathcal{Y}|} \sum_{y \in \mathcal{Y}} y \right)^2 \right] \\ &= \frac{1 + (N/|\mathcal{X}| + |\mathcal{Y}|)^{-1}}{1 + |\mathcal{Y}|^{-1}} \left[ \left( \frac{1}{|\mathcal{Y}|} \sum_{y \in \mathcal{Y}} y^2 \right) - \left( \frac{1}{|\mathcal{Y}|} \sum_{y \in \mathcal{Y}} y \right)^2 \right].\end{aligned}\quad (3.47)$$

Since all possible values of  $x$  are equally probable and the conditional probability  $P_{y|x}$  is uniform and independent of  $x$ , the risk simply becomes the variance of the set  $\mathcal{Y}$  scaled by a factor dependent on  $|\mathcal{Y}|$  and on  $N/|\mathcal{X}|$ . Without training data ( $N = 0$ ), the scaling is unity; as  $N/|\mathcal{X}| \rightarrow \infty$ , the scaling factor is  $(1 + |\mathcal{Y}|^{-1})^{-1}$ .

To visualize the performance, use the explicit sets  $\mathcal{Y}$  and  $\mathcal{X}$  defined earlier. The conditional variance becomes

$$\Sigma_{y|x} = \frac{|\mathcal{Y}|^2 - 1}{12|\mathcal{Y}|^2} = \frac{1 - |\mathcal{Y}|^{-2}}{12} \quad (3.48)$$

and the minimum risk is expressed as

$$\begin{aligned}\mathcal{R}^* &= \frac{(1 - |\mathcal{Y}|^{-1})(1 + (N/|\mathcal{X}| + |\mathcal{Y}|)^{-1})}{12} \\ &= \left( \frac{|\mathcal{Y}|}{N/|\mathcal{X}| + |\mathcal{Y}|} \right) \frac{1 - |\mathcal{Y}|^{-2}}{12} + \left( \frac{N/|\mathcal{X}|}{N/|\mathcal{X}| + |\mathcal{Y}|} \right) \frac{1 - |\mathcal{Y}|^{-1}}{12}.\end{aligned}\quad (3.49)$$

Interestingly, the minimum squared-error for the uniform prior can be represented as a convex combination of two separate risk values with weighting factors dependent on  $|\mathcal{Y}|$  and  $N/|\mathcal{X}|$ . Thus for a uniform prior, the risk depends on the number of elements in  $\mathcal{Y}$  and the number of training samples “per element of  $\mathcal{X}$ ”. Note the relationship of these weighting factors to those of the conditional PMF  $P_{y|x,D}$ , which depend on  $\alpha'(x)$  and on  $N'(x; D)$ . For the uniform prior,  $\alpha'(x) = |\mathcal{Y}|$  and  $E_D [N'(D)] = N/|\mathcal{X}|$ .

The first risk is the conditional variance  $\Sigma_{y|x}$  - this is intuitively satisfying as the corresponding weight becomes unity when  $N = 0$ . The second risk is the squared-error with infinite training data. Note that the reduction of the risk between these two extreme cases is modest, and that the attenuating factor increases towards unity for applications with more possible outcomes. Figure 3.18 illustrates the difference between these cases.

PGR: additional figures for uniform case?

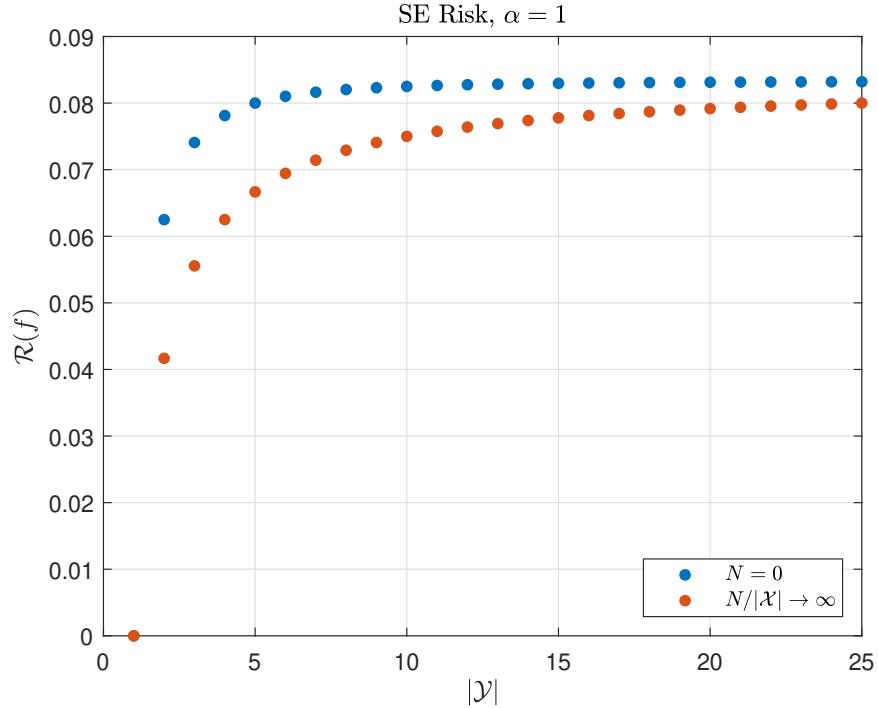


Figure 3.18: Minimum SE Risk, Uniform Prior, zero and infinite training data

### 3.3.1.3 Conditional Squared-Error for a Dirichlet-based Estimator

Having derived the optimal estimator based on a Dirichlet model prior, it is important to consider the conditional risk  $\mathcal{R}_\Theta(f^*; \theta)$  and analyze how different prior parametrizations  $\alpha$  influence the squared-error for different models  $\theta$ . Starting from the conditional squared-error risk (2.28) and substituting the Bayesian estimator (2.31), the formula simplifies to

$$\begin{aligned}\mathcal{R}_\Theta(f^*; \theta) &= \mathcal{R}_\Theta^*(\theta) + E_{x,D|\theta} \left[ (f^*(x; D) - f_\Theta(x; \theta))^2 \right] \\ &= E_{x|\theta} [\Sigma_{y|x,\theta}] + E_{x,D|\theta} \left[ (\mu_{y|x,D} - \mu_{y|x,\theta})^2 \right].\end{aligned}\quad (3.50)$$

Defining the excess conditional risk  $\mathcal{R}_{\Theta,\text{ex}}(f; \theta) \equiv \mathcal{R}_\Theta(f; \theta) - \mathcal{R}_\Theta^*(\theta)$ , the second term above is  $\mathcal{R}_{\Theta,\text{ex}}(f^*; \theta) = E_{x,D|\theta} \left[ (\mu_{y|x,D} - \mu_{y|x,\theta})^2 \right]$ , the average squared bias between the Bayesian predictive mean and the true predictive mean.

PGR: general above? move before dir? simplify form in terms of Ptheta?

PGR: define excess clairvoyant/bayesian risk up front before dir?

Evaluation of the excess risk for an estimator based on the Dirichlet prior will be performed using the sufficient statistic  $\bar{n}$  in place of the training set  $D$ . Using the random

process  $\Delta(x, \bar{n}, \theta) \equiv P_{y|x,\bar{n}} - P_{y|x,\theta} \in \mathbb{R}^{\mathcal{Y}}$  introduced in 3.2, the term is expressed as

$$\begin{aligned}\mathcal{R}_{\Theta,\text{ex}}(f^*; \theta) &= E_{x,D|\theta} \left[ (\mu_{y|x,D} - \mu_{y|x,\theta})^2 \right] \\ &= \sum_{y \in \mathcal{Y}} y \sum_{y' \in \mathcal{Y}} y' E_{x,\bar{n}|\theta} \left[ \Delta(y; x, \bar{n}, \theta) \Delta(y'; x, \bar{n}, \theta) \right] \\ &= \sum_{y \in \mathcal{Y}} y \sum_{y' \in \mathcal{Y}} y' E_{x,n'|\theta} \left[ \mathcal{E}(y, y'; x, n', \theta) \right] \\ &= E_{x|\theta} \left[ \sum_{y|x,\theta} E_{n'(x)|\theta'(x)} \left[ \frac{n'(x)}{(\alpha'(x) + n'(x))^2} \right] \right] \\ &\quad + E_{x|\theta} \left[ (\mu_{y|x} - \mu_{y|x,\theta})^2 E_{n'(x)|\theta'(x)} \left[ \left( \frac{\alpha'(x)}{\alpha'(x) + n'(x)} \right)^2 \right] \right],\end{aligned}\tag{3.51}$$

where the function  $\mathcal{E}$  is defined in (3.39).

The excess conditional risk can thus be represented as the conditional expectation (with respect to  $P_{x|\theta}$ ) of a sum of two functions of  $x$ . The first function measures the additional variance beyond that of the clairvoyant estimator (i.e. the clairvoyant squared-error); like the clairvoyant risk, it depends on  $\Sigma_{y|x,\theta}$ , the conditional variance of the clairvoyant estimate for a given observation of  $x$ . The second function is dependent on the squared bias between the clairvoyant estimate  $\mu_{y|x,\theta}$  and the data-independent estimate  $\mu_{y|x}$ . This term alone is influenced by the data-independent Bayes predictive distribution  $P_{y|x} = \alpha(\cdot, x)/\alpha'(x)$ .

The two second-order (in terms of  $y$ ) terms are scaled by factors dependent on the prior concentrations  $\alpha'(x)$  and on  $\theta'(x)$  and  $N$  via conditional expectations with respect to  $n'(x)$ . Note that by the aggregation property of multinomial distributions, the random variable  $n'(x)|\theta'(x) \sim Bi(N, \theta'(x))$ . Closed-forms have not been found for the function expectations of binomial random variables above.

PGR: binomial citations?

It is instructional to consider the trends in the conditional squared-error risk (3.51) for different volumes of training data  $N$  and for different selections of  $\alpha$ .

First consider how the excess risk changes with the training volume  $N$ . For  $N = 0$ , it is evident that the excess risk is  $\mathcal{R}_{\Theta,\text{ex}}(f^*; \theta) \rightarrow E_{x|\theta} \left[ (\mu_{y|x} - \mu_{y|x,\theta})^2 \right]$ , the expected squared bias between the clairvoyant and data-independent estimators. As  $N$  tends to infinity, the binomial distribution controlling the scaling factors concentrates at  $n'(x) \approx N\theta'(x)$ ; as such, the two expectations of interest tend to zero and thus  $\mathcal{R}_{\Theta,\text{ex}}(f^*; \theta) \rightarrow 0$ . This desirable

property of the estimator is a consequence of the full support of the Dirichlet prior, ensuring that the model posterior concentrates at the empirical PMF.

Another interesting point regarding the dependency of the excess conditional risk on  $N$  is that, depending on the learner parameterization, there may be a local maximum. Consider the trivial case of  $|\mathcal{X}| = 1$  - treating  $N$  as a real number, there would be a maximum at

$$N \equiv \alpha'(x) \left( 1 - 2\alpha'(x) \frac{(\mu_{y|x} - \mu_{y|x,\theta})^2}{\Sigma_{y|x,\theta}} \right). \quad (3.52)$$

Note that as the squared-difference between the prior mean and true mean increases, the maximizing value decreases (even below zero). Thus, the worse the prior estimate, the more likely the excess squared-error will decrease monotonically with  $N$ . Conversely, if the prior estimate is accurate, a local maximum may occur and additional training data may (temporarily) compromise the estimator performance. Also consider the effect of prior concentration; informative priors with sufficiently high  $\alpha'(x)$  will not have the local maxima.

The excess risk at this potentially non-integral value would be

$$\mathcal{R}_{\Theta,\text{ex}}(f^*; \theta) \rightarrow E_{x|\theta} \left[ \frac{\frac{1}{\alpha'(x)} - \frac{(\mu_{y|x} - \mu_{y|x,\theta})^2}{\Sigma_{y|x,\theta}}}{4 \left( 1 - \frac{(\mu_{y|x} - \mu_{y|x,\theta})^2}{\Sigma_{y|x,\theta}} \right)} \Sigma_{y|x,\theta} \right]. \quad (3.53)$$

PGR: better form above???

Figures 3.19 and 3.20 exemplify the excess conditional squared-error as a function of  $N$  for estimators based on Dirichlet priors of varying concentration  $\alpha'(x)$ . The former shows local maxima for an unbiased estimator; note that higher concentration results in superior performance. The latter uses biased estimators and as such, learners based on low concentration achieve lower risk.

Next consider the effects of the Dirichlet prior parameters. The analysis will interpret the Dirichlet parameters as the conditional prior distributions  $\alpha(\cdot, x)/\alpha'(x)$  and their concentrations  $\alpha'(x)$ .

First consider the conditional prior PMF's  $\alpha(\cdot, x)/\alpha'(x)$ ; as shown, they manifest themselves in the risk through the squared estimator bias. It is clear that regardless of how the values  $\alpha'(x)$  are chosen, the best selections for these conditional priors must have first moments matching those of the corresponding clarivoyant predictive distributions  $P_{y|x,\theta}$  for

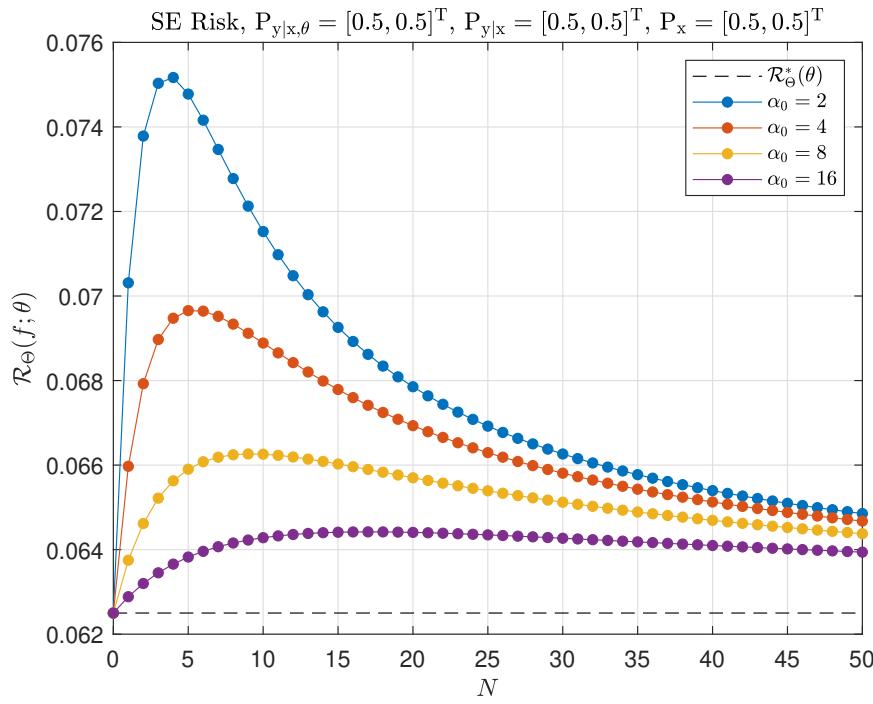


Figure 3.19: Conditional SE Risk versus  $N$ , unbiased Dirichlet estimators of varying concentration

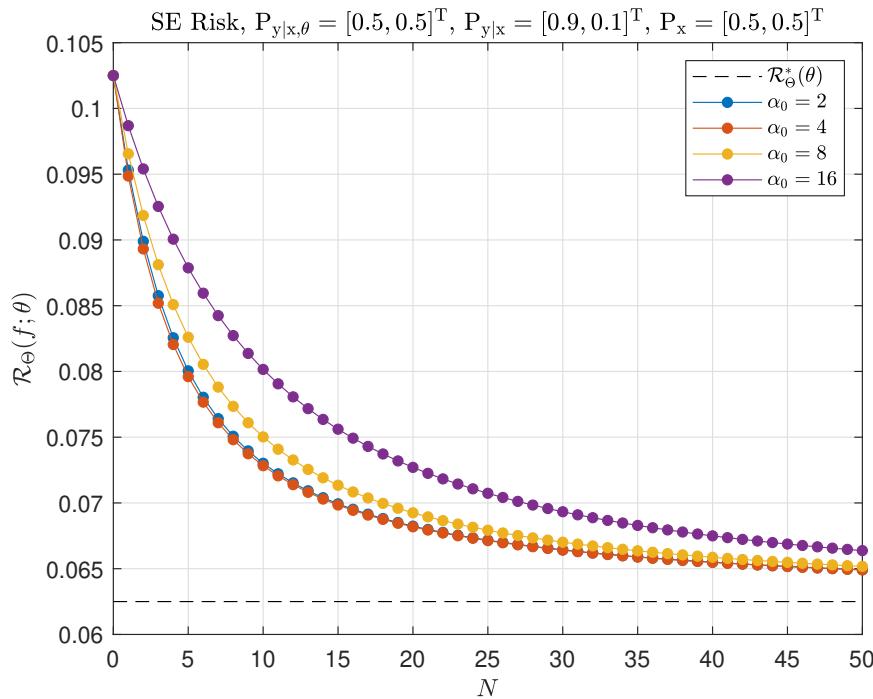


Figure 3.20: Conditional SE Risk versus  $N$ , biased Dirichlet estimators of varying concentration

each  $x \in \mathcal{X}$ . Such estimators are unbiased; as a result, the excess conditional risk is thus equivalent to the first term in (3.51), measuring additional variance due to model uncertainty.

The other user-selected Dirichlet parameters  $\alpha'(x)$  are the concentration parameters for the corresponding conditional distributions; they control important bias/variance trade-offs via the two scaling factors in (3.51). First, consider the asymptotic trends.

Consider how the excess risk tends as the priors become maximally concentrated. As the parameters  $\alpha'(x) \rightarrow \infty$ , the excess risk tends to  $\mathcal{R}_{\Theta,\text{ex}}(f^*; \theta) \rightarrow E_{x|\theta} \left[ (\mu_{y|x} - \mu_{y|x,\theta})^2 \right]$ , the expected conditional squared-error between the means of the Bayesian predictive PMF and the clairvoyant predictive PMF. This is intuitive given that the estimator tends toward a data-independent solution; analogous to the discussion in Section 3.2, the estimator may be biased, but will have no variance due to the training data statistics.

Conversely, if concentrations  $\alpha'(x) \rightarrow 0$  are chosen, the Bayesian estimate tends to the empirical mean, independent of  $\alpha(\cdot, x)/\alpha'(x)$ , and the excess risk tends to

$$\begin{aligned} \mathcal{R}_{\Theta,\text{ex}}(f^*; \theta) &\rightarrow E_{x|\theta} \left[ \sum_{n=1}^N \binom{N}{n} \theta'(x)^n (1 - \theta'(x))^{N-n} \frac{1}{n} \right] \\ &\quad + E_{x|\theta} \left[ (1 - \theta'(x))^N (\mu_{y|x} - \mu_{y|x,\theta})^2 \right]. \end{aligned}$$

Note that the first term's scaling factor is proportionate to the first inverse moment of a positive binomial random variable [17]. The second term's scaling factor tends to  $P_{n'(x)|\theta'(x)}(0|\theta'(x))$ , the probability that no training samples are observed matching the value  $x$ . As  $N$  increases, this term tends to zero, the risk due to the prior estimate bias decreases, and the excess risk becomes a function of  $\theta$  only.

Of further interest are the values  $\alpha'(x)$  that minimize the excess squared-error for a given prior conditional distribution  $\alpha(\cdot, x)/\alpha'(x)$ . With the asymptotic values of the excess risk known, all that remains is to determine any local minima. Since the excess risk is a sum of  $|\mathcal{X}|$  terms of identical form, each dependent on their own concentration  $\alpha'(x)$ , only one component needs to be minimized.

PGR: add the derivative details below???

Calculating the first derivative with respect to  $\alpha'(x)$ , it can be shown that for  $N > 0$  and

$\theta'(x) > 0$ , only one stationary point exists, at

$$\alpha'(x) \equiv \frac{\Sigma_{y|x,\theta}}{(\mu_{y|x} - \mu_{y|x,\theta})^2}. \quad (3.54)$$

Calculation of the second derivative confirms that this value is a local minimum. Furthermore, the excess risk evaluated at these values is

$$\mathcal{R}_{\Theta,\text{ex}}(f^*; \theta) = E_{x|\theta} \left[ E_{n'(x)|\theta'(x)} \left[ \frac{1}{n'(x)\Sigma_{y|x,\theta}^{-1} + (\mu_{y|x} - \mu_{y|x,\theta})^{-2}} \right] \right], \quad (3.55)$$

which can be easily shown to be less than both the asymptotic values for  $\alpha'(x) \rightarrow 0$  and  $\alpha'(x) \rightarrow \infty$ . Thus the concentration values (3.54) provide the minimum excess risk for the given prior conditional distributions.

Note that the minimizing concentration values  $\alpha'(x)$  are inversely proportional to the squared-bias of the prior conditional mean. This is sensible; the better the match between the true and prior predictive distributions, the more confidence should be expressed. Also, low concentrations are preferable when the model has low conditional variance; these models can be quickly identified with learners prioritizing the empirical PMF estimate over the prior estimate. Additionally, note that these values  $\alpha'(x)$  do not depend on the training volume  $N$ .

Figures 3.21 and 3.22 show how the excess conditional squared-error trends as a function of the Dirichlet learner concentration. Note that the latter is based on a biased prior estimate and thus the optimal Dirichlet concentration value is lower.

PGR: plot captions, alpha zero or x???

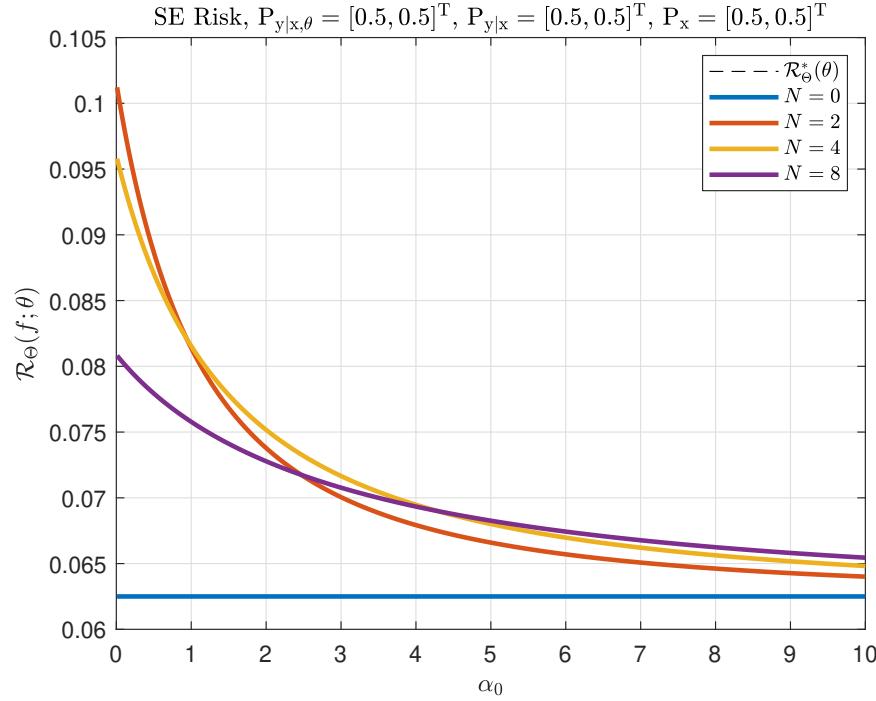


Figure 3.21: Conditional SE Risk versus  $\alpha'(x)$ , unbiased Dirichlet estimator using varying training set volumes

PGR: newpage

### 3.3.2 Classification: the 0-1 Loss

This section derives 0-1 loss classifiers based on the Dirichlet prior distribution and assesses their performance.

#### 3.3.2.1 Optimal Hypothesis: Conditional Maximum *a posteriori*

PGR: decision region figures??

PGR: weighted conditional majority decision

To determine the optimal learning function, the 0-1 loss from Equation (2.34) is substi-

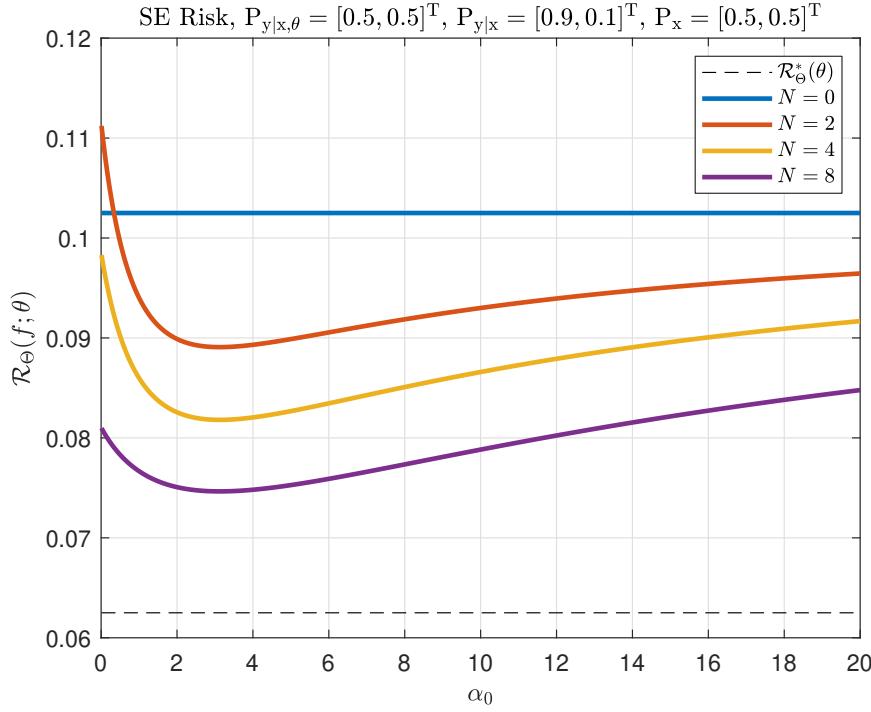


Figure 3.22: Conditional SE Risk versus  $\alpha'(x)$ , biased Dirichlet estimator using varying training set volumes

tuted into Equation (3.41) and Equation (2.8) to find

$$\begin{aligned} f^*(x; D) &= \arg \max_{y \in \mathcal{Y}} P_{y|x,D}(y|x, D) \\ &= \arg \max_{y \in \mathcal{Y}} \frac{\alpha(y, x) + \bar{N}(y, x; D)}{\alpha'(x) + N'(x; D)} \\ &= \arg \max_{y \in \mathcal{Y}} (\alpha(y, x) + \bar{N}(y, x; D)) . \end{aligned} \quad (3.56)$$

Using the Dirichlet prior, different classes are “scored” by counting the number of training samples with a value of  $X_n$  matching that of  $x$  and combining with the prior parameters  $\alpha(\cdot, x)$ .

**Uniform Prior** When the uniform prior is used, the Bayes classifier simplifies to

$$f^*(x; D) = \arg \max_{y \in \mathcal{Y}} \bar{N}(y, x; D) , \quad (3.57)$$

a conditional majority decision which chooses the class from  $\mathcal{Y}$  most often represented among training set samples  $D$  with a matching input value  $x$ . This is intuitive, as the model PDF parameter  $\alpha$  imparts no confidence as to which classes may be most likely.

### 3.3.2.2 Minimum Risk: Probability of Error

PGR: GENERATE NON-SIM FIGS!!!!!!

PGR: DIR SIM FIGS COMMENTED!!!

PGR: no closed-forms found???

Evaluating the minimum risk (2.40) using the distributions derived from the Dirichlet prior, the Bayes minimum probability of error is

$$\begin{aligned}\mathcal{R}^* &= 1 - E_{x,D} \left[ \max_{y \in \mathcal{Y}} P_{y|x,D}(y|x, D) \right] \\ &= 1 - E_{x,\bar{n}} \left[ \frac{\max_{y \in \mathcal{Y}} (\alpha(y, x) + \bar{n}(y, x))}{\alpha'(x) + n'(x)} \right] \\ &= 1 - \sum_{x \in \mathcal{X}} \frac{E_{\bar{n}} \left[ \max_{y \in \mathcal{Y}} (\alpha(y, x) + \bar{n}(y, x)) \right]}{\alpha_0 + N}.\end{aligned}\quad (3.58)$$

Figures 3.23 and 3.24 plot the minimum Bayes probability of error against training data volume  $N$  and prior concentration  $\alpha_0$ , respectively. Note that for  $N = 0$ , the Bayes risk is  $\mathcal{R}^* = 1 - \sum_{x \in \mathcal{X}} \frac{\max_{y \in \mathcal{Y}} \alpha(y, x)}{\alpha_0}$ . Additionally, consider the risk for maximal/minimal values of the Dirichlet concentration. For  $\alpha_0 \rightarrow 0$  (and  $N > 1$ ), the risk is  $\mathcal{R}^* = 0$ ; conversely, for  $\alpha_0 \rightarrow \infty$ , the risk tends to  $\mathcal{R}^* \rightarrow 1 - \sum_{x \in \mathcal{X}} \frac{\max_{y \in \mathcal{Y}} \alpha(y, x)}{\alpha_0}$ . These trends can be visualized in Figures 3.25 and 3.26.

PGR: risk for  $N \rightarrow \infty$ ?

PGR: missing info for Dir gen graphics? fixed y given x conditional alpha???

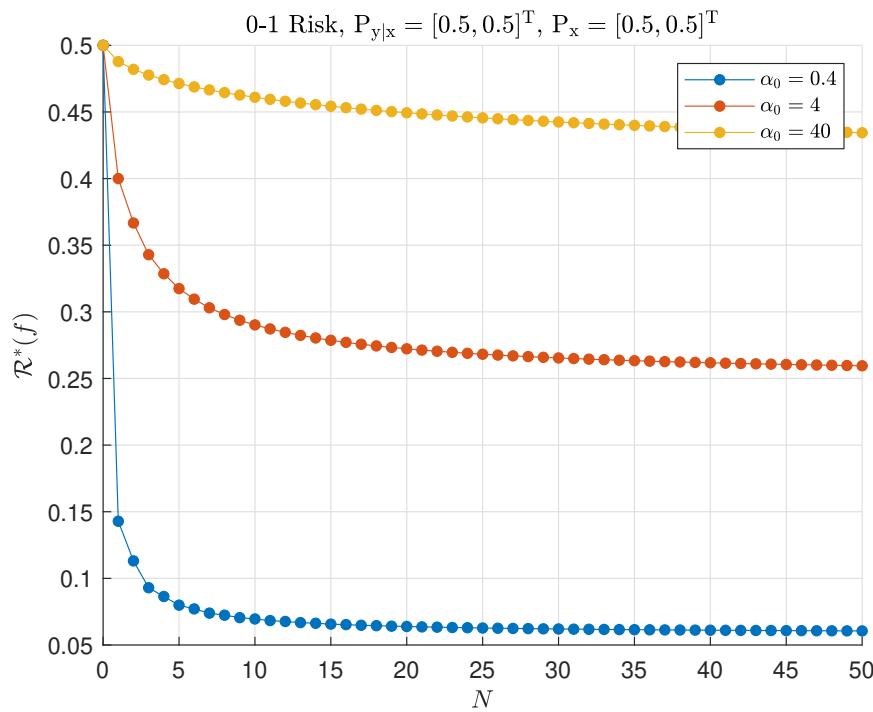
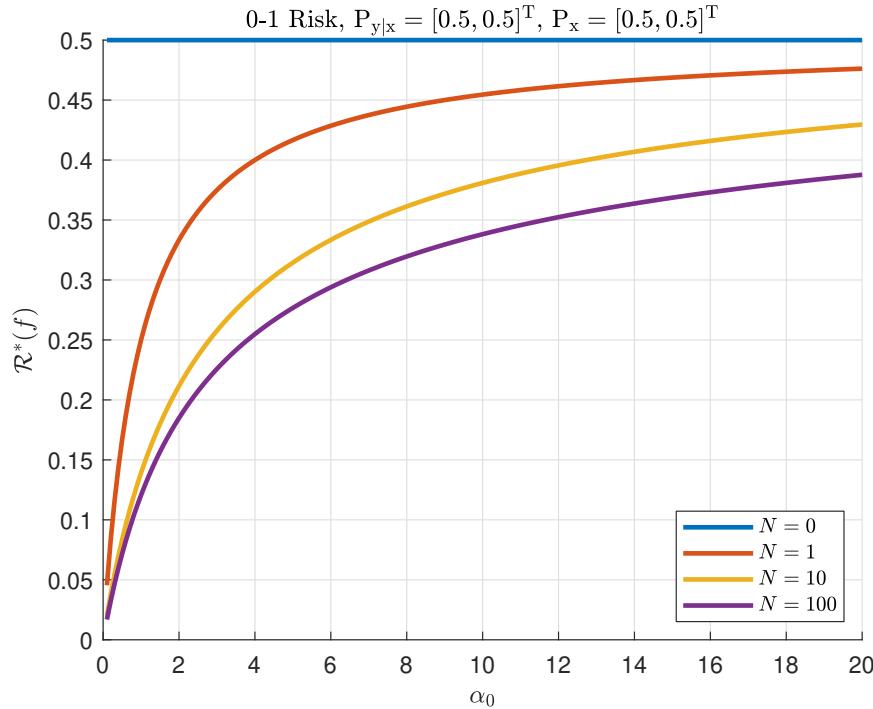
PGR: comment on simulation!

**Uniform Prior** PGR: COMPUTATIONAL COMPLEXITY savings for risk formula?

PGR: Can uniform minimal risk be approximated as a function of  $M_y$  and  $M_x/N$ , as is for SE loss???

PGR: use Mcal not binom!

PGR: add nmax CDF fig!

Figure 3.23: Minimum 0-1 Risk for different training data volumes  $N$ Figure 3.24: Minimum 0-1 Risk for different prior concentrations  $\alpha_0$

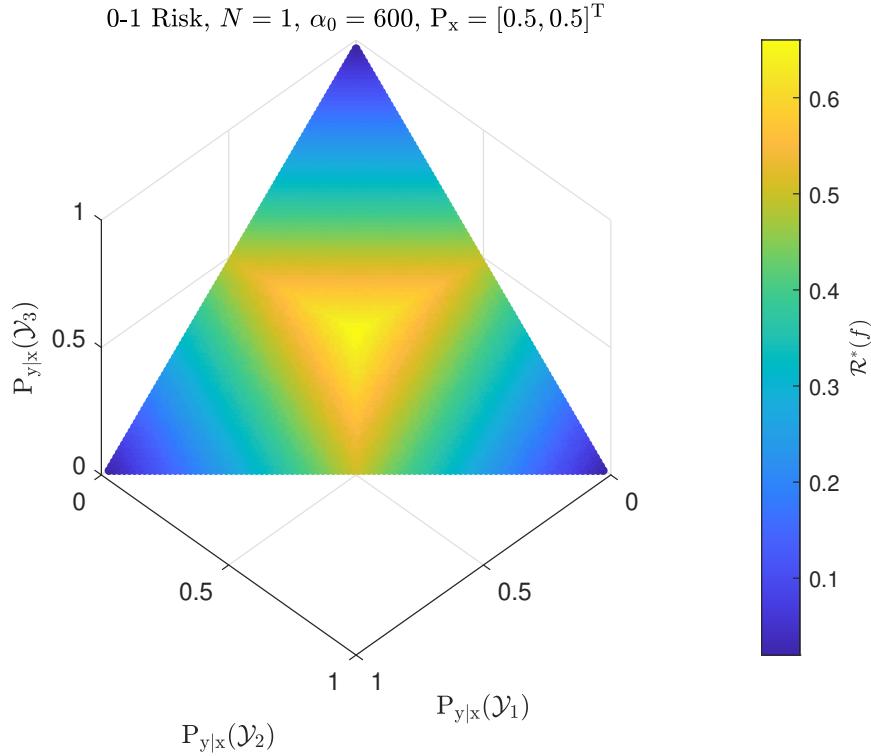


Figure 3.25: Minimum 0-1 Risk for different prior means  $P_{y|x}$

Using the uniform prior, the minimum Bayes 0-1 risk is

$$\begin{aligned}
 R^* &= 1 - E_{x,D} \left[ \max_{y \in \mathcal{Y}} P_{y|x,D}(y|x, D) \right] \\
 &= 1 - \sum_{x \in \mathcal{X}} \frac{E_{\bar{n}} \left[ \max_{y \in \mathcal{Y}} \bar{n}(y, x) \right] + 1}{|\mathcal{Y}| |\mathcal{X}| + N} \\
 &= 1 - \frac{1 + |\mathcal{X}|^{-1} \sum_{x \in \mathcal{X}} E_{\bar{n}} \left[ \max_{y \in \mathcal{Y}} \bar{n}(y, x) \right]}{|\mathcal{Y}| + N/|\mathcal{X}|}.
 \end{aligned} \tag{3.59}$$

The expectation operates on the maximum value from a subset of a uniform Dirichlet-Multinomial random process. Via the Dirichlet-Multinomial aggregation property [9], a consequence of the the uniform PMF  $P_{\bar{n}}$  is that the individual segments  $\bar{n}(\cdot, x)$  are identically distributed; thus, the expectation will be same for every value  $x$ .

To evaluate this expectation, new random variables  $\bar{n}_{\max}(x) \equiv \max_{y \in \mathcal{Y}} \bar{n}(y, x)$  are introduced and characterized by their identical PMF. To this end, the probability of the event  $P(\bar{n}_{\max}(x) \geq n) = P(\cup_{y \in \mathcal{Y}} \{\bar{n}(y, x) \geq n\})$  will be determined. As the distribution of  $\bar{n}$  is uniform, the event probability is proportionate to the cardinality of the set

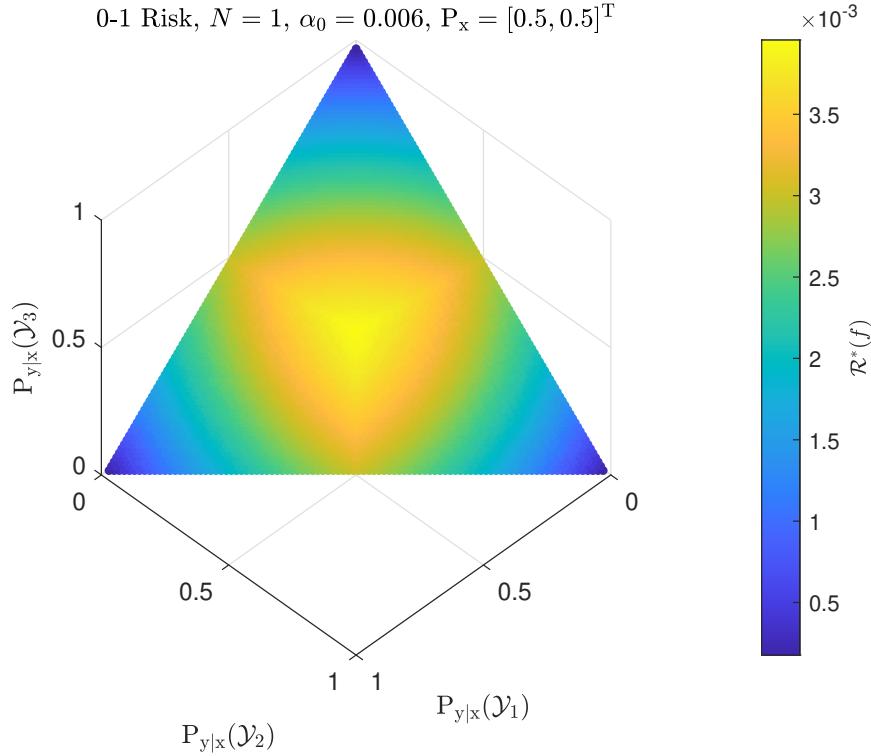


Figure 3.26: Minimum 0-1 Risk for different prior means  $P_{y|x}$

$\cup_{y \in \mathcal{Y}} \{\bar{n} : \bar{n}(y, x) \geq n\}$ . Using the inclusion-exclusion principle [5], the cardinality is represented as

$$\begin{aligned}
 & \left| \cup_{y \in \mathcal{Y}} \{\bar{n} : \bar{n}(y, x) \geq n\} \right| \tag{3.60} \\
 &= \begin{cases} \binom{N+|\mathcal{Y}||\mathcal{X}|-1}{|\mathcal{Y}||\mathcal{X}|-1} & \text{if } n < 0, \\ \sum_{m=1}^{|\mathcal{Y}|} \binom{|\mathcal{Y}|}{m} (-1)^{m-1} \binom{N-mn+|\mathcal{Y}||\mathcal{X}|-1}{|\mathcal{Y}||\mathcal{X}|-1} H\left(\left\lfloor \frac{N}{m} \right\rfloor - n\right) & \text{if } 0 \leq n \leq N, \\ 0 & \text{if } n > N, \end{cases}
 \end{aligned}$$

where  $H : \mathbb{Z} \mapsto \{0, 1\}$  is the discrete Heaviside step function. For  $n < 0$ , the cardinality is equivalent to  $|\bar{\mathcal{N}}|$ .

For  $0 \leq n < N$ , the cardinality is an alternating binomial summation where the  $m^{\text{th}}$  term accounts for the different intersections of  $m$  of the  $|\mathcal{Y}|$  individual sets  $\{\bar{n} : \bar{n}(y, x) \geq n\}$ . Observe that the cardinality of the intersections is only dependent on the number of contributing sets  $m$  and not on which sets intersect. Furthermore, note the dependency of the intersection cardinalities on the argument  $n$ . The step function contributes such

that if  $n > \lfloor \frac{N}{m} \rfloor$ , only up to  $m - 1$  individual sets will intersect. The binomial coefficient  $\mathcal{M}(\{N - mn, |\mathcal{Y}||\mathcal{X}| - 1\})$  provides the intersection cardinality for a given  $m$ ; note the similarity to the cardinality  $|\bar{\mathcal{N}}|$  - the only difference is the number of points characterizing the  $|\mathcal{Y}||\mathcal{X}| - 1$  dimensional region.

The probability of interest can thus be expressed as

$$\begin{aligned} P(\bar{n}_{\max}(x) \geq n) &= \binom{N + |\mathcal{Y}||\mathcal{X}| - 1}{|\mathcal{Y}||\mathcal{X}| - 1}^{-1} |\cup_{y \in \mathcal{Y}} \{\bar{n} : \bar{n}(y, x) \geq n\}| \\ &= \begin{cases} 1 & \text{if } n < 0, \\ \sum_{m=1}^{|\mathcal{Y}|} \binom{|\mathcal{Y}|}{m} (-1)^{m-1} \prod_{l=1}^{|\mathcal{Y}||\mathcal{X}|-1} \left(1 - \frac{mn}{N+l}\right) H\left(\lfloor \frac{N}{m} \rfloor - n\right) & \text{if } 0 \leq n \leq N, \\ 0 & \text{if } n > N. \end{cases} \end{aligned} \quad (3.61)$$

PGR: use Mcal op?

PGR: Heaviside reference?

As the PMF of  $\bar{n}_{\max}(x)$  has support on  $n \in [0, \dots, N]$ , the expectation over  $\bar{n}$  is evaluated as

$$\begin{aligned} E_{\bar{n}} [\bar{n}_{\max}(x)] &= \sum_{n=0}^N n \left( P(\bar{n}_{\max}(x) \geq n) - P(\bar{n}_{\max}(x) \geq n+1) \right) \\ &= -1 + \sum_{n=0}^N P(\bar{n}_{\max}(x) \geq n) \\ &= -1 + \sum_{m=1}^{|\mathcal{Y}|} \binom{|\mathcal{Y}|}{m} (-1)^{m-1} \sum_{n=0}^{\lfloor \frac{N}{m} \rfloor} \prod_{l=1}^{|\mathcal{Y}||\mathcal{X}|-1} \left(1 - \frac{mn}{N+l}\right) \end{aligned} \quad (3.62)$$

and the minimum 0-1 risk is

$$\mathcal{R}^* = 1 - \frac{\sum_{m=1}^{|\mathcal{Y}|} \binom{|\mathcal{Y}|}{m} (-1)^{m-1} \sum_{n=0}^{\lfloor \frac{N}{m} \rfloor} \prod_{l=1}^{|\mathcal{Y}||\mathcal{X}|-1} \left(1 - \frac{mn}{N+l}\right)}{|\mathcal{Y}| + N/|\mathcal{X}|}. \quad (3.63)$$

It is instructional to express the risk for minimal and maximal volumes of training data.

Using the binomial summation identity

$$\sum_{m=0}^M \binom{M}{m} (-1)^m g(m) = 0, \quad (3.64)$$

where  $g$  is a polynomial function of degree less than  $M$  [8], it can be shown that for  $N = 0$ , the minimum risk is  $\mathcal{R}^* = 1 - |\mathcal{Y}|^{-1}$ . This is sensible, as the classes are equiprobable with  $P_y = |\mathcal{Y}|^{-1}$ .

PGR: use ruiz citation for identity?

PGR: find irreducible risk explicitly from theta?

To find the risk for  $N \rightarrow \infty$ , note that

$$\begin{aligned} & \lim_{N \rightarrow \infty} (|\mathcal{Y}| + N/|\mathcal{X}|)^{-1} \sum_{n=0}^{\lfloor \frac{N}{m} \rfloor} \prod_{l=1}^{|\mathcal{Y}||\mathcal{X}|-1} \left(1 - \frac{mn}{N+l}\right) \\ &= \lim_{N/m \rightarrow \infty} \frac{|\mathcal{X}|}{m} \sum_{n=0}^{\lfloor \frac{N}{m} \rfloor} \left(1 - \frac{mn}{N}\right)^{|\mathcal{Y}||\mathcal{X}|-1} \frac{m}{N} \\ &= \frac{|\mathcal{X}|}{m} \int_0^1 (1-t)^{|\mathcal{Y}||\mathcal{X}|-1} dt \\ &= \frac{1}{m|\mathcal{Y}|}. \end{aligned} \quad (3.65)$$

The irreducible 0-1 risk for the uniform prior tends toward

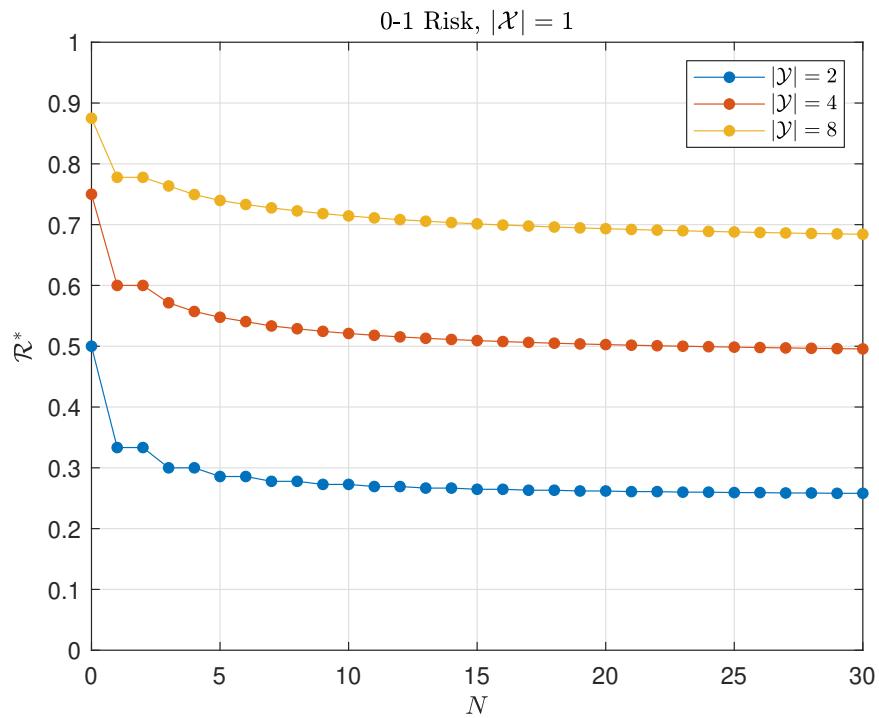
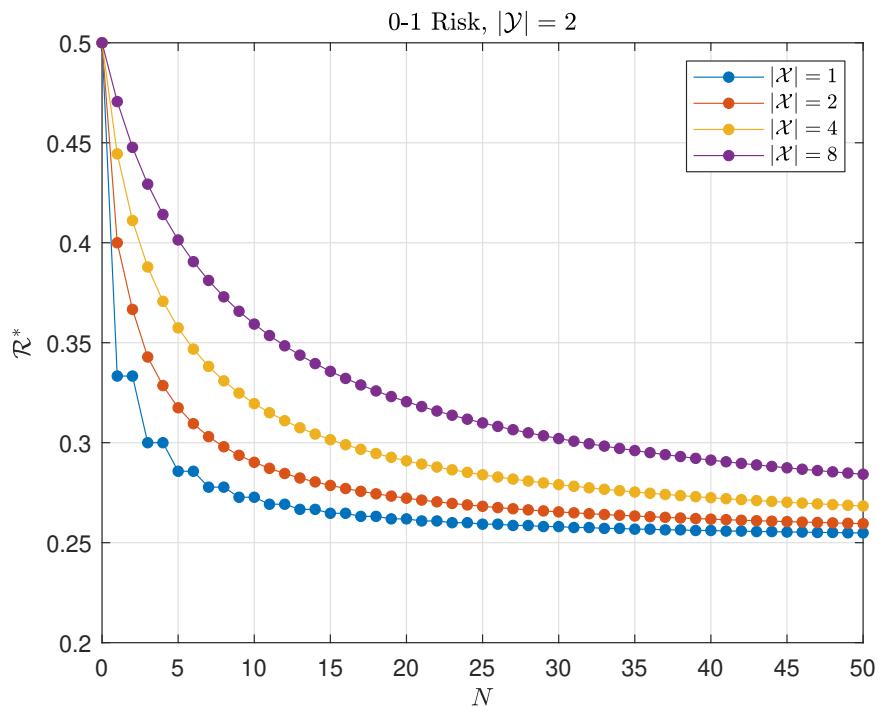
$$\begin{aligned} \mathcal{R}^* &\rightarrow 1 - |\mathcal{Y}|^{-1} \sum_{m=1}^{|\mathcal{Y}|} \binom{|\mathcal{Y}|}{m} (-1)^{m-1} m^{-1} \\ &= 1 - |\mathcal{Y}|^{-1} \sum_{m=1}^{|\mathcal{Y}|} m^{-1}, \end{aligned} \quad (3.66)$$

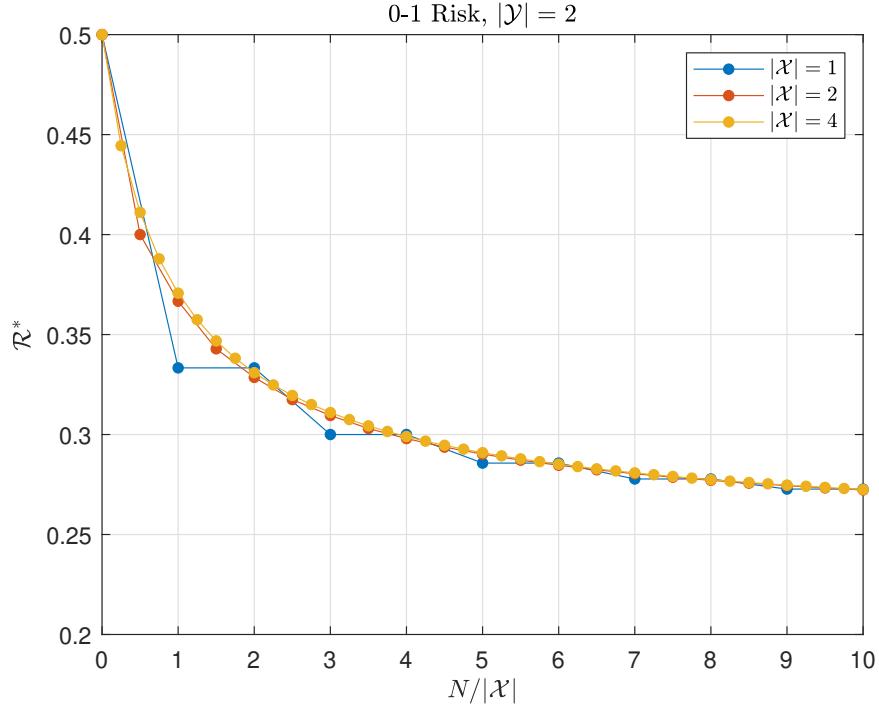
providing a lower bound for the achievable 0-1 Bayes risk. The above formulation has made use of the alternating summation identity from [15] to display the risk with a form including the  $|\mathcal{Y}|^{\text{th}}$  harmonic number  $H_{|\mathcal{Y}|} \equiv \sum_{m=1}^{|\mathcal{Y}|} m^{-1}$ . Observe that the irreducible risk does not depend on the cardinality  $|\mathcal{X}|$ .

PGR: harmonic reference?

Figure 3.27 demonstrates how the minimum 0-1 risk decreases with training volume  $N$ ; observe that the risk is more severe for sequences corresponding to higher  $|\mathcal{Y}|$ . It is sensible that the probability of error should increase when more classes have to be considered. Figure 3.28 illustrates the minimum risk with multiple sequences for different cardinalities  $|\mathcal{X}|$ . Note that risk increases with  $|\mathcal{X}|$ . Considering  $E_D[N'(D)] = \mu_{n'} = N/|\mathcal{X}|$ , this should be intuitive - each conditional empirical distribution  $\bar{N}(\cdot, x; D)/N'(x; D)$  is forced to approximate  $\tilde{\theta}(x)$  with less data.

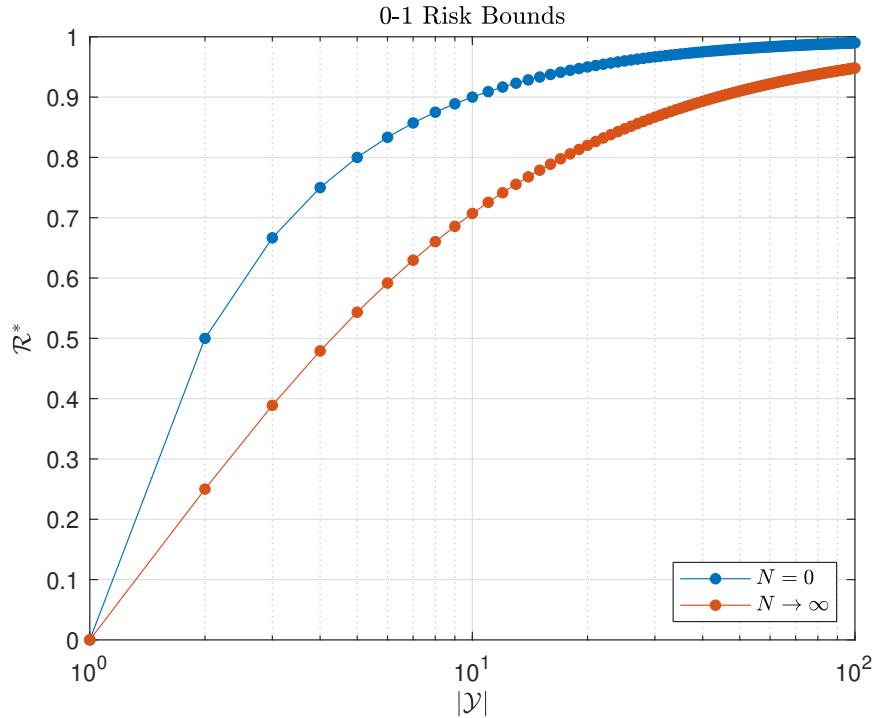
Further insight into how  $|\mathcal{X}|$  affects the risk can be acquired by plotting the risk as a function of  $N/|\mathcal{X}|$ . In Figure 3.29, it is shown that the minimal risk can be approximated

Figure 3.27: Minimum 0-1 Risk vs training set volume  $N$ Figure 3.28: Minimum 0-1 Risk vs training set volume  $N$

Figure 3.29: Minimum 0-1 Risk vs  $N/|\mathcal{X}|$ 

by a function dependent only on  $N/|\mathcal{X}|$ ; of the series plotted, only the series for  $|\mathcal{X}| = 1$  shows notable non-negligible from the others.

It is also useful to graph the  $N = 0$  and  $N \rightarrow \infty$  minimum risk as a function of  $|\mathcal{Y}|$ ; both formulas are independent of  $|\mathcal{X}|$ . Figure 3.30 displays these bounds; note the margin in the probability of error between the optimal  $N = 0$  and  $N \rightarrow \infty$  classifiers. For  $|\mathcal{Y}| = 2$  binary classification, both sequences are at their minimum and infinite training data provides a reduction in expected probability of error from 0.5 to 0.25. As  $|\mathcal{Y}|$  increases, the classification risk for both the  $N = 0$  and  $N \rightarrow \infty$  cases tend to unity and the error reduction for  $N \rightarrow \infty$  decreases.

Figure 3.30: Minimum 0-1 Risk vs  $|\mathcal{Y}|$ 

PGR: newpage

### 3.3.2.3 Conditional Probability of Error for a Dirichlet-based Classifier

PGR: INCOMPLETE

PGR: comment on alpha0 versus alphax simplification

Substituting the optimal Dirichlet-based classifier into the formula for the conditional probability of error 2.35, the risk is

$$\mathcal{R}_\theta(f; \theta) = 1 - \sum_{x \in \mathcal{X}} \theta'(x) E_{\bar{n}|\theta} \left[ \tilde{\theta} \left( \arg \max_{y \in \mathcal{Y}} (\bar{n}(y, x) + \alpha(y, x)); x \right) \right]. \quad (3.67)$$

Figures 3.31 and 3.32 show how the conditional risk trends for classifiers based on well-matched and poorly-matched informative Dirichlet priors, respectively. Note that the well-matched prior does better with higher prior concentrations  $\alpha_0$ ; this is reflective of the fact that the maximizing arguments  $y \in \mathcal{Y}$  of both the true model  $\tilde{\theta}(x)$  and the prior mean  $\alpha(\cdot, x)/\alpha'(x)$  are the same.

Also, it is important to consider how a given classifier performs for varying models  $\theta$ .

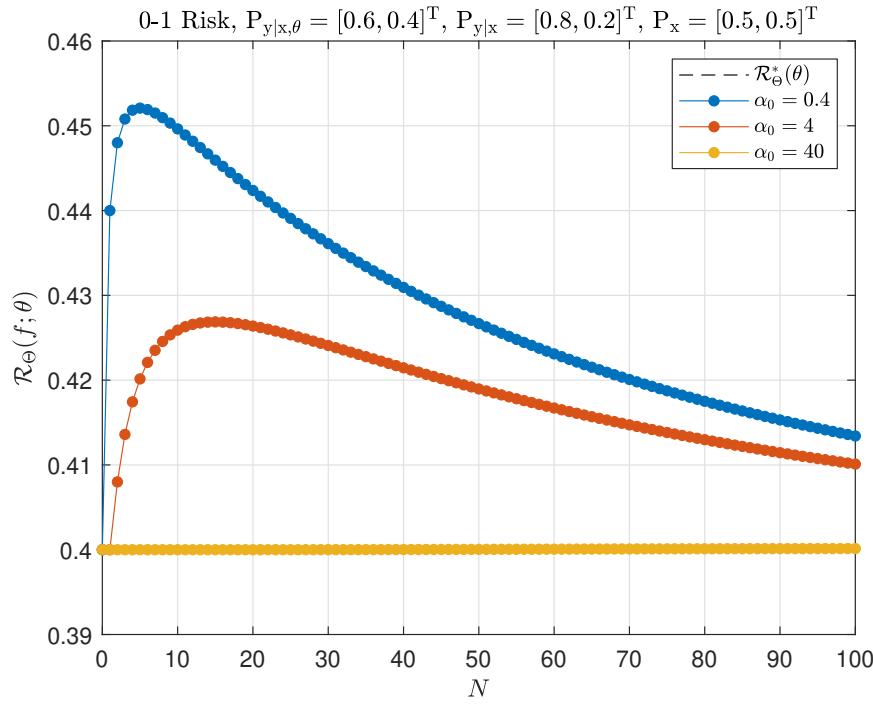


Figure 3.31: Excess conditional probability of error, well-matched informative Dirichlet-based classifier

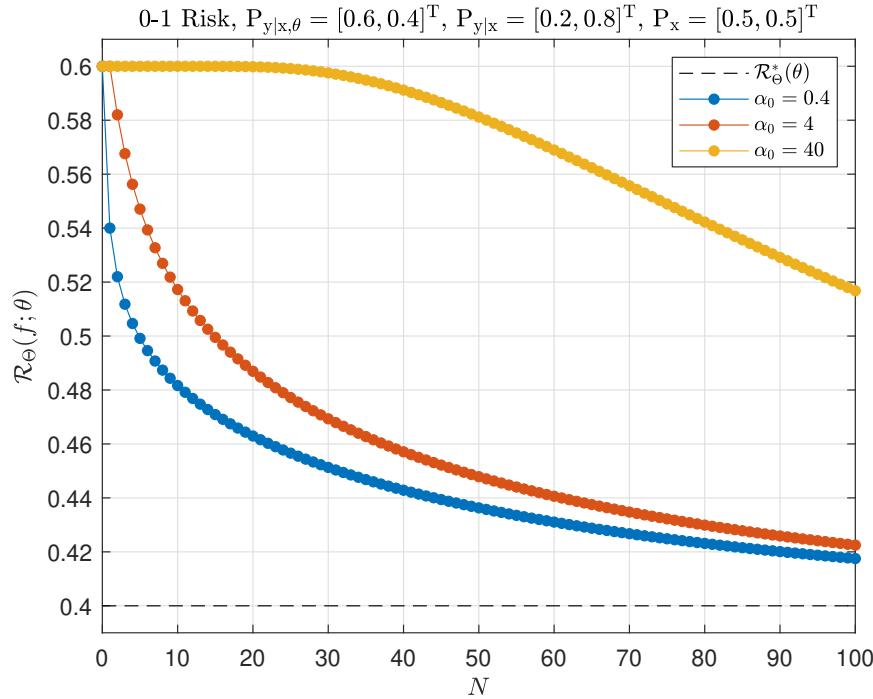


Figure 3.32: Excess conditional probability of error, poorly-matched informative Dirichlet-based classifier

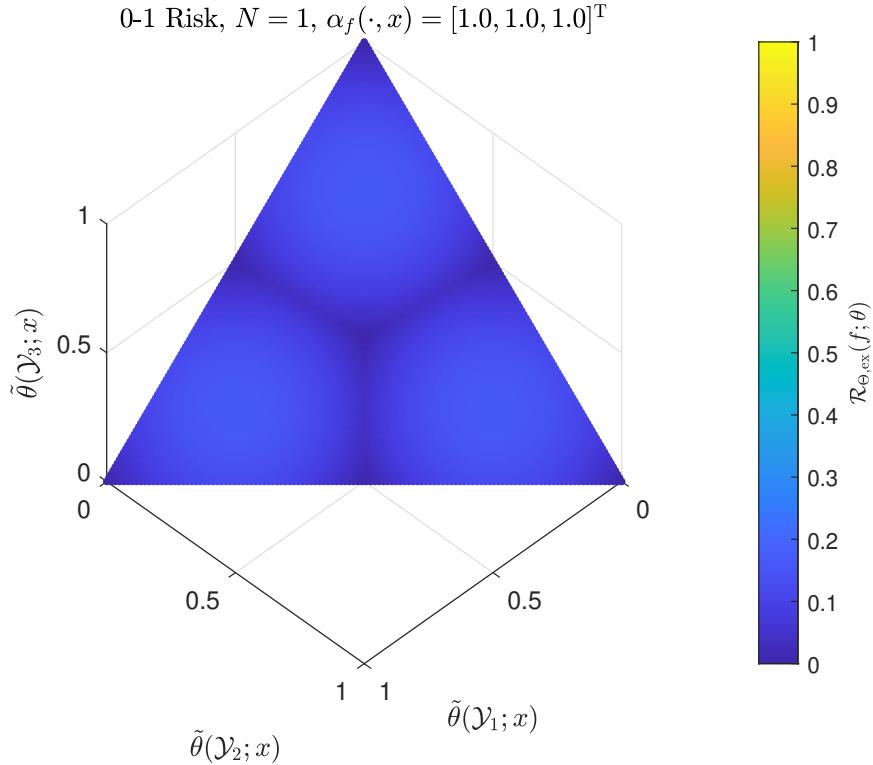


Figure 3.33: Excess conditional probability of error, conditional majority decision

Figures 3.33 and 3.34 demonstrate the excess conditional probability of error achieved by the conditional majority decision (based on a non-informative Dirichlet prior) and by a classifier derived from an informative Dirichlet prior, respectively. Note that while the former has fewer models for which the error is critically high, the latter has more models for which the clairvoyant risk  $\mathcal{R}_{\Theta}^*(\theta)$  is achieved. This a fundamental trade-off between Bayesian learners based on non-informative versus informative priors.

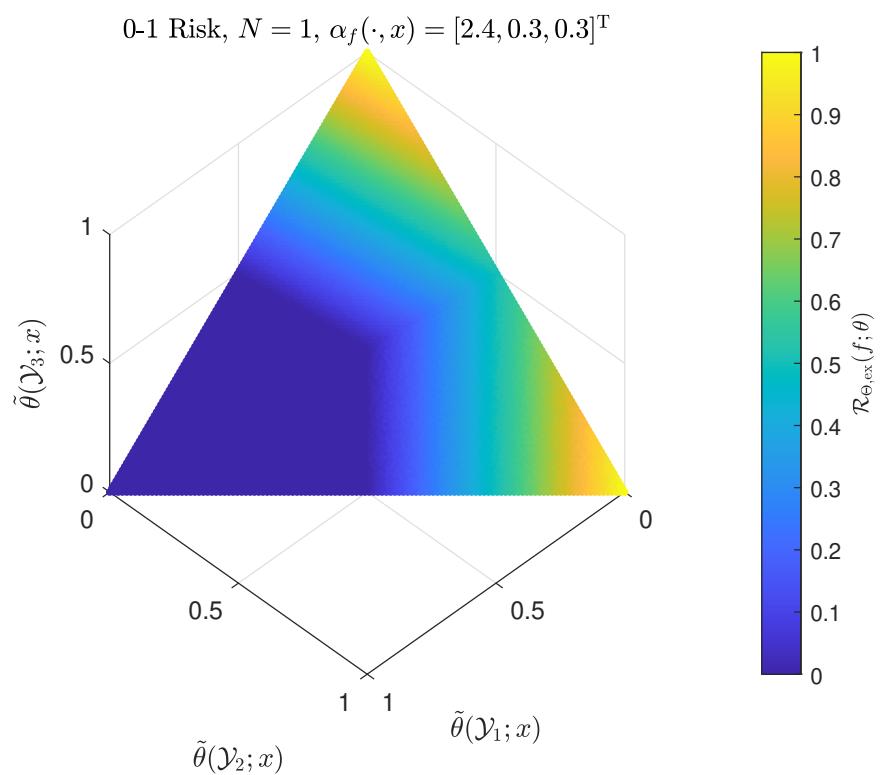


Figure 3.34: Excess conditional probability of error, informative Dirichlet-based classifier

# Chapter 4

## Extention to Infinite-Dimensional Spaces - Countably Infinite

PGR: MERGE with finite chapter???

### 4.1 Intro

This chapter extends previous results for applications where the space  $\mathcal{Y}$  is countably infinite, that is  $|\mathcal{Y}| = \aleph_0$ . Specifically, the model prior distribution will be characterized by a discrete-domain Dirichlet process.

### 4.2 Basic Model

#### 4.2.1 Probability Distributions

PGR: ???

##### 4.2.1.1 Model PDF, $p(\theta)$

PGR: Valid model representation? Marginals instead?

$$p(\theta) = \beta(\alpha)^{-1} \prod_{y \in \mathcal{Y}} \theta(y)^{\alpha(y)-1}, \quad (4.1)$$

$$\beta(\alpha) = \frac{\prod_{y \in \mathcal{Y}} \Gamma(\alpha(y))}{\Gamma\left(\sum_{y \in \mathcal{Y}} \alpha(y)\right)}. \quad (4.2)$$

The first and second joint moments of the model are

$$\mu_\theta(y) = E_\theta [\theta(y)] = \frac{\alpha(y)}{\alpha_0} \quad (4.3)$$

and

$$E_\theta [\theta(y)\theta(y')] = \frac{\alpha(y)\alpha(y') + \alpha(y)\delta[y, y']}{\alpha_0(\alpha_0 + 1)}. \quad (4.4)$$

#### 4.2.1.2 Training Data PMF, P(D)

$$P(D | \theta) = \prod_{y \in \mathcal{Y}} \theta(y)^{\bar{N}(y; D)}. \quad (4.5)$$

$$P(\bar{n} | \theta) = \mathcal{M}(\bar{n}) \prod_{y \in \mathcal{Y}} \theta(y)^{\bar{n}(y)}, \quad (4.6)$$

Thus,

$$P(\bar{n}) = \mathcal{M}(\bar{n}) \beta(\alpha)^{-1} \beta(\alpha + \bar{n}). \quad (4.7)$$

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

$$E_{\bar{n}} [\bar{n}(y)] = N \frac{\alpha(y)}{\alpha_0} \quad (4.8)$$

and

$$E_{\bar{n}} [\bar{n}(y)\bar{n}(y')] = \frac{N}{\alpha_0(\alpha_0 + 1)} ((\alpha_0 + N)\alpha(y)\delta[y, y'] + (N - 1)\alpha(y)\alpha(y')). \quad (4.9)$$

Also,

$$P(D) = \beta(\alpha)^{-1} \beta(\alpha + \bar{N}(D)). \quad (4.10)$$

#### 4.2.1.3 Output conditional PMF, $P(y|D)$

$$\begin{aligned} p(\theta|D) &= \frac{P(D|\theta)p(\theta)}{P(D)} \\ &= \beta (\alpha + \bar{N}(D))^{-1} \prod_{y \in \mathcal{Y}} \theta(y)^{\alpha(y) + \bar{N}(y; D) - 1} \end{aligned} \quad (4.11)$$

$$p(\theta|\bar{n}) = \beta (\alpha + \bar{n})^{-1} \prod_{y \in \mathcal{Y}} \theta(y)^{\alpha(y) + \bar{n}(y) - 1}, \quad (4.12)$$

The PMF of interest is

$$\begin{aligned} P(y|D) &= E_{\theta|D} [\theta(y)] \\ &= \frac{\alpha(y) + \bar{N}(y; D)}{\alpha_0 + N} \\ &= \left( \frac{\alpha_0}{\alpha_0 + N} \right) \frac{\alpha(y)}{\alpha_0} + \left( \frac{N}{\alpha_0 + N} \right) \frac{\bar{N}(y; D)}{N} \end{aligned} \quad (4.13)$$

$$\begin{aligned} P(y|\bar{n}) &= E_{\theta|\bar{n}} [\theta(y)|\bar{n}] \\ &= \frac{\alpha(y) + \bar{n}(y)}{\alpha_0 + N} \\ &= \left( \frac{\alpha_0}{\alpha_0 + N} \right) \frac{\alpha(y)}{\alpha_0} + \left( \frac{N}{\alpha_0 + N} \right) \frac{\bar{n}(y)}{N} \end{aligned} \quad (4.14)$$

### 4.3 Application to Common Loss Functions

PGR: Definitely regression. Classification sensible for countably infinite?

PGR: Results identical to finite Dirichlet

$$\begin{aligned} E_{y|D} [\mathcal{L}(h, y)] &= \sum_{y \in \mathcal{Y}} \mathcal{L}(h, y) P_{y|D}(y|D) \\ &= \frac{\sum_{y \in \mathcal{Y}} \alpha(y) \mathcal{L}(h, y) + \sum_{y \in \mathcal{Y}} \bar{N}(y; D) \mathcal{L}(h, y)}{\alpha_0 + N} \\ &= \frac{\sum_{y \in \mathcal{Y}} \alpha(y) \mathcal{L}(h, y) + \sum_{n=1}^N \mathcal{L}(h, D_n)}{\alpha_0 + N} \\ &= \left( \frac{\alpha_0}{\alpha_0 + N} \right) \sum_{y \in \mathcal{Y}} \mathcal{L}(h, y) \frac{\alpha(y)}{\alpha_0} + \left( \frac{N}{\alpha_0 + N} \right) N^{-1} \sum_{n=1}^N \mathcal{L}(h, D_n). \end{aligned} \quad (4.15)$$

## 4.4 General Model

Extension to output and input spaces  $\mathcal{Y}$  and  $\mathcal{X}$  can have an infinite number of elements.

## 4.5 Applications: General Model

PGR: Definitely regression. Classification sensible for countably infinite?

PGR: Results identical to finite Dirichlet

# Chapter 5

## Extention to Infinite-Dimensional Spaces - Uncountably Infinite

PGR: SPECIFY EUCLIDEAN/HILBERT??

PGR: account for impulsive alpha?

### 5.1 Intro

This chapter extends further to the case where  $\mathcal{Y}$  is a Euclidean space and the model  $\theta$  is a continuous-domain Dirichlet process.

### 5.2 Basic Model

#### 5.2.1 Probability Distributions

PGR: ???

##### 5.2.1.1 Model $\theta$ Characterization

The model is now a continuous-domain Dirichlet process  $\theta \sim DP(\alpha)$ . The concentration parameter is  $\alpha_0 \equiv \int_{y \in \mathcal{Y}} \alpha(y) dy$ . By definition, for any partition of the set  $\mathcal{Y}$ ,  $\{\dots, \mathcal{S}(z), \dots\}$ ,  $z \in \mathcal{Z}$ , we can generate a Dirichlet random process  $\phi(z) \equiv \int_{\mathcal{S}(z)} \theta(y) dy$  with parameterizing

function  $\lambda(z) \equiv \int_{\mathcal{S}(z)} \alpha(y) dy$ . This is commonly referred to as the aggregation property. The PDF for the aggregation is thus

$$P_\phi(\phi) = \beta(\lambda)^{-1} \prod_{z \in \mathcal{Z}} \phi(z)^{\lambda(z)-1}. \quad (5.1)$$

As detailed in Appendix A.4, the first and second moments of the Dirichlet process  $\Theta$  are

$$\mu_\Theta = \frac{\alpha}{\alpha_0} \quad (5.2)$$

and

$$E_\Theta [\Theta(y)\Theta(y')] = \frac{\alpha(y)\alpha(y') + \alpha(y)\delta(y-y')}{\alpha_0(\alpha_0+1)}. \quad (5.3)$$

Again,  $p_y = p_{D_n} = \mu_\Theta$  for all  $n$ .

### 5.2.1.2 Output conditional PDF, $p_{y|D}$

In Appendix A.5, it was shown that if the model  $\Theta \sim DP(\alpha)$  is a Dirichlet process, then the model conditioned on the training data  $D$  is also a Dirichlet process with parameterizing function  $\alpha + \bar{N}(D)$ , where  $\bar{N}(y; D) = \sum_{n=1}^N \delta(y - D_n)$  generalizes for a continuous domain. Thus, the conditional PDF of interest can be formulated as

$$\begin{aligned} p_{y|D} &= \mu_{\Theta|D} \\ &= \frac{\alpha + \bar{N}(D)}{\alpha_0 + N} = \frac{\alpha + \sum_{n=1}^N \delta(\cdot - D_n)}{\alpha_0 + N} \\ &= \left( \frac{\alpha_0}{\alpha_0 + N} \right) \frac{\alpha}{\alpha_0} + \left( \frac{N}{\alpha_0 + N} \right) \frac{\sum_{n=1}^N \delta(\cdot - D_n)}{N} \\ &= \left( \frac{\alpha_0}{\alpha_0 + N} \right) p_y + \left( \frac{N}{\alpha_0 + N} \right) \frac{\sum_{n=1}^N \delta(\cdot - D_n)}{N}. \end{aligned} \quad (5.4)$$

With the generalization to a Euclidean space  $\mathcal{Y}$ , the training data dependent component of the PDF is formulated with Dirac delta functions.

### 5.2.1.3 Training Data PDF, $p(D)$

To represent the training data distribution, note that the Dirichlet process conditional model also provides

$$p_{D_{n+1}|D_n, \dots, D_1} = \frac{\alpha + \sum_{i=1}^n \delta(\cdot - D_i)}{\alpha_0 + N} \quad (5.5)$$

and thus the complete PDF is

$$\begin{aligned} p_D(D) &= p_{D_1}(D_1) \prod_{n=2}^N p_{D_n|D_{n-1}, \dots, D_1}(D_n | D_{n-1}, \dots, D_1) \\ &= \frac{\alpha(D_1)}{\alpha_0} \prod_{n=2}^N \frac{\alpha(D_n) + \sum_{i=1}^{n-1} \delta(D_n - D_i)}{\alpha_0 + n - 1}. \end{aligned} \quad (5.6)$$

Additionally, note that since  $p_{D_n|\theta} = \theta$  is independent of sample index  $n$ , the PDF does not vary when the input arguments are permuted. Furthermore, all marginal distributions of  $D$  will have the same form, regardless of which training samples  $D_n$  are used.

Using these properties, the first and second joint moments of  $D$  are found to be

$$\begin{aligned} \mu_{D_n} &= \int_y y p_{D_n}(y) dy = \int_y y E_\theta [P_{D_n|\theta}(y)] dy \\ &= \int_y y \mu_\theta(y) dy \\ &= \int_y y \frac{\alpha(y)}{\alpha_0} dy \equiv \mu_y \end{aligned} \quad (5.7)$$

and

$$\begin{aligned} E_D [D_n^2] &= \int_y y^2 p_{D_n}(y) dy = \int_y y^2 E_\theta [P_{D_n|\theta}(y)] dy \\ &= \int_y y^2 \mu_\theta(y) dy \\ &= \int_y y^2 \frac{\alpha(y)}{\alpha_0} dy = E_y[y^2], \end{aligned} \quad (5.8)$$

$$\begin{aligned} E_D [D_n D_{n'}] &= \int_y \int_y yy' p_{D_n, D_{n'}}(y, y') dy dy' \\ &= \int_y \int_y yy' E_\theta [p_{D_n|\theta}(y) p_{D_{n'}|\theta}(y')] dy dy' \\ &= \int_y \int_y yy' E_\theta [\theta(y) \theta(y')] dy dy' \\ &= \int_y \int_y yy' \frac{\alpha(y) \alpha(y') + \alpha(y) \delta(y - y')}{\alpha_0(\alpha_0 + 1)} dy dy' \\ &= \frac{\alpha_0 \mu_y^2 + E_y[y^2]}{\alpha_0 + 1}. \end{aligned} \quad (5.9)$$

Combining,

$$E_D [D_n D_{n'}] = E_y[y^2] - (1 - \delta[n, n']) \frac{\alpha_0}{\alpha_0 + 1} \Sigma_y. \quad (5.10)$$

PGR: move proofs to appendix???

PGR PGR: Dirichlet-Multinomial Process perspective???

Define the Dirichlet-Multinomial process  $\bar{n} \equiv \bar{N}(D)$ . The mean and correlation functions below are found in Appendix A.6. The mean function is

$$\mu_{\bar{n}} = N \frac{\alpha}{\alpha_0} \quad (5.11)$$

and the correlation function is

$$E_{\bar{n}} [\bar{n}(y)\bar{n}(y')] = \frac{N}{\alpha_0(\alpha_0 + 1)} [(N - 1)\alpha(y)\alpha(y') + (\alpha_0 + N)\alpha(y)\delta(y - y')] . \quad (5.12)$$

## 5.3 Application to Common Loss Functions

PGR: Discuss continuous notation

$$\begin{aligned} E_{y|D} [\mathcal{L}(h, y)] &= \int_{\mathcal{Y}} \mathcal{L}(h, y) p_{y|D}(y|D) dy \\ &= \frac{\int_{\mathcal{Y}} \alpha(y) \mathcal{L}(h, y) dy + \int_{\mathcal{Y}} \sum_{n=1}^N \delta(y - D_n) \mathcal{L}(h, y) dy}{\alpha_0 + N} \\ &= \frac{\int_{\mathcal{Y}} \alpha(y) \mathcal{L}(h, y) dy + \sum_{n=1}^N \mathcal{L}(h, D_n)}{\alpha_0 + N} \\ &= \left( \frac{\alpha_0}{\alpha_0 + N} \right) \int_{\mathcal{Y}} \frac{\alpha(y)}{\alpha_0} \mathcal{L}(h, y) dy + \left( \frac{N}{\alpha_0 + N} \right) N^{-1} \sum_{n=1}^N \mathcal{L}(h, D_n) \\ &= \left( \frac{\alpha_0}{\alpha_0 + N} \right) E_y [\mathcal{L}(h, y)] + \left( \frac{N}{\alpha_0 + N} \right) N^{-1} \sum_{n=1}^N \mathcal{L}(h, D_n) \end{aligned} \quad (5.13)$$

### 5.3.1 Regression: the Squared-Error Loss

$$\mathcal{L}(h, y) = (h - y)^2 . \quad (5.14)$$

Now we choose for the regression function to map to  $\mathcal{H} = \mathcal{Y} = \mathbb{R}$ .

$$\begin{aligned}
 \mathcal{R}(f) &= E_{\theta} \left[ E_{D|\theta} \left[ E_{y|\theta} \left[ (f(D) - y)^2 \right] \right] \right] \\
 &= E_{\theta} \left[ E_{y|\theta} [(y - \mu_{y|\theta})^2] \right] + E_{\theta} \left[ E_{D|\theta} \left[ (f(D) - \mu_{y|\theta})^2 \right] \right] \\
 &= E_{\theta} [\Sigma_{y|\theta}] + E_{\theta} \left[ E_{D|\theta} \left[ (f(D) - \mu_{y|\theta})^2 \right] \right]
 \end{aligned} \tag{5.15}$$

### 5.3.1.1 Optimal Learner

The optimal function is again the expected value of the output conditional PMF,

$$\begin{aligned}
 f^*(D) &= \arg \min_{h \in \mathbb{R}} E_{y|D} [(h - y)^2] \\
 &= \mu_{y|D} = E_{\theta|D} [\mu_{y|\theta}] \\
 &= \left( \frac{\alpha_0}{\alpha_0 + N} \right) \int_{\mathcal{Y}} \frac{\alpha(y)}{\alpha_0} y dy + \left( \frac{N}{\alpha_0 + N} \right) \frac{1}{N} \sum_{n=1}^N D_n \\
 &= \left( \frac{\alpha_0}{\alpha_0 + N} \right) \mu_y + \left( \frac{N}{\alpha_0 + N} \right) \frac{1}{N} \sum_{n=1}^N D_n \\
 &= \left( \frac{\alpha_0}{\alpha_0 + N} \right) \mu_y + \left( \frac{N}{\alpha_0 + N} \right) \int_{\mathcal{Y}} y \frac{\bar{N}(y; D)}{N} dy .
 \end{aligned} \tag{5.16}$$

### 5.3.1.2 Minimum Risk

$$\begin{aligned}
 \mathcal{R}^* &= E_D [\Sigma_{y|D}] \\
 &= E_{\theta} [\Sigma_{y|\theta}] + E_D \left[ E_{\theta|D} \left[ (\mu_{y|\theta} - E_{\theta|D} [\mu_{y|\theta}])^2 \right] \right] \\
 &= E_{\theta} [\Sigma_{y|\theta}] + E_D [C_{\theta|D} [\mu_{y|\theta}]] \\
 &= E_{\bar{n}} [\Sigma_{y|\bar{n}}] .
 \end{aligned} \tag{5.17}$$

The conditional variance is expanded as

$$\Sigma_{y|D} = E_{y|D}[y^2] - \mu_{y|D}^2 \tag{5.18}$$

and the expectations of the two terms are evaluated separately.

$$E_D [E_{y|D}[y^2]] = E_y[y^2] \tag{5.19}$$

$$\begin{aligned}
 & E_D [\mu_y^2|_D] \\
 &= \frac{\alpha_0^2 \mu_y^2 + 2\alpha_0 \mu_y \sum_{n=1}^N \mu_{D_n} + \sum_{n=1}^N \sum_{n'=1}^N E_D [D_n D_{n'}]}{(\alpha_0 + N)^2} \\
 &= \frac{\alpha_0^2 \mu_y^2 + 2\alpha_0 N \mu_y^2 + N^2 E_y[y^2] - N(N-1)\alpha_0(\alpha_0+1)^{-1}\Sigma_y}{(\alpha_0 + N)^2} \\
 &= \mu_y^2 + \frac{N}{(\alpha_0 + 1)(\alpha_0 + N)} \Sigma_y
 \end{aligned} \tag{5.20}$$

PGR: DMP PERSPECTIVE???

$$\begin{aligned}
 E_D [\mu_y^2|_D] &= E_{\bar{n}} [\mu_y^2|_{\bar{n}}] \\
 &= \frac{\alpha_0^2 \mu_y^2 + 2\alpha_0 \mu_y \int_{\mathcal{Y}} y \mu_{\bar{n}}(y) dy + \int_{\mathcal{Y}} \int_{\mathcal{Y}} yy' E_{\bar{n}} [\bar{n}(y)\bar{n}(y')] dy dy'}{(\alpha_0 + N)^2} \\
 &= \frac{\alpha_0^2 \mu_y^2 + 2\alpha_0 N \mu_y^2 + N(N-1)\alpha_0(\alpha_0+1)^{-1}\mu_y^2 + N(\alpha_0+N)(\alpha_0+1)^{-1} E_y[y^2]}{(\alpha_0 + N)^2} \\
 &= \frac{\alpha_0(\alpha_0+N+1)\mu_y^2 + N E_y[y^2]}{(\alpha_0 + 1)(\alpha_0 + N)}
 \end{aligned} \tag{5.21}$$

PGR: nicer algebra with DMP!

PGR: DMP

The minimal risk is again

$$\begin{aligned}
 \mathcal{R}^* &= \left(1 - \frac{N}{(\alpha_0 + 1)(\alpha_0 + N)}\right) \Sigma_y \\
 &= \frac{\alpha_0(\alpha_0+N+1)}{(\alpha_0 + 1)(\alpha_0 + N)} \Sigma_y \\
 &= \frac{1 + (\alpha_0 + N)^{-1}}{1 + \alpha_0^{-1}} \Sigma_y.
 \end{aligned} \tag{5.22}$$

As for Dirichlet distributions for functions with countable domains, the minimal risk is dependent on the model only through the concentration parameter  $\alpha_0$  and the variance of the expected distribution  $P_y = \mu_\theta$ .

## 5.4 General Model

PGR: change dirac deltas to kronecker in fractions, no divide by zero???

### 5.4.1 Model Extension

This section adds the input space  $\mathcal{X}$ , now considered to be uncountably infinite; that is  $|\mathcal{X}| \geq \aleph_1$ . The model distribution is a Dirichlet process over the space  $\mathcal{Y} \times \mathcal{X}$ .

### 5.4.2 General Probability Distributions

PGR

#### 5.4.2.1 Model $\theta$ Characterization

The model is characterized by a Dirichlet process  $\theta \sim DP(\alpha)$  with parameter function  $\alpha : \mathcal{Y} \times \mathcal{X} \mapsto \mathbb{R}^+$ . The concentration parameter generalizes to  $\alpha_0 = \int_{\mathcal{Y}} \int_{\mathcal{X}} \alpha(y, x) dx dy$ . Using the aggregation property, any partition of the set  $\mathcal{Y} \times \mathcal{X}$ ,  $\{\dots, \mathcal{S}(z), \dots\}$ ,  $z \in \mathcal{Z}$  is a Dirichlet random process  $\phi(z) = \int \int_{\mathcal{S}(z)} \theta(y, x) dx dy$  with parameterizing function  $\lambda(z) = \int \int_{\mathcal{S}(z)} \alpha(y, x) dx dy$ . The PDF for the aggregation is

$$p_{\phi}(\phi) = \beta(\lambda)^{-1} \prod_{z \in \mathcal{Z}} \phi(z)^{\lambda(z)-1}. \quad (5.23)$$

The expected value of a Dirichlet process generalizes to  $DP(\alpha)$

$$\mu_{\theta} = \frac{\alpha}{\alpha_0} \quad (5.24)$$

and the correlation function is

$$E_{\theta} [\theta(y, x)\theta(y', x')] = \frac{\alpha(y, x)\alpha(y', x') + \alpha(y, x)\delta(y - y')\delta(x - x')}{\alpha_0(\alpha_0 + 1)}. \quad (5.25)$$

#### 5.4.2.2 Output conditional PDF, $p_{y|x,D}$

The properties of the Dirichlet distribution proven in Appendix A.5 generalize, such that the model conditioned on the training data  $D$  is Dirichlet with parameterizing function  $\alpha + \bar{N}(D)$ , where  $\bar{N}(y, x; D) = \sum_{n=1}^N \delta(y - Y_n)(x - X_n)$ . Recall that  $D_n = (Y_n, X_n)$  with  $Y \in \mathcal{Y}^N$  and  $X \in \mathcal{X}^N$ .

The PDF  $p_{y,x|D}$  is thus

$$\begin{aligned} p_{y,x|D} &= \mu_{\theta|D} \\ &= \frac{\alpha + \bar{N}(D)}{\alpha_0 + N} \\ &= \frac{\alpha + \sum_{n=1}^N \delta(\cdot - Y_n) \delta(\cdot - X_n)}{\alpha_0 + N} \\ &= \left( \frac{\alpha_0}{\alpha_0 + N} \right) \frac{\alpha}{\alpha_0} + \left( \frac{N}{\alpha_0 + N} \right) \frac{\sum_{n=1}^N \delta(\cdot - Y_n) \delta(\cdot - X_n)}{N} \end{aligned} \quad (5.26)$$

and the conditional PDF of interest is

$$\begin{aligned} p_{y|x,D} &= \frac{\alpha(\cdot, x) + \bar{N}(\cdot, x; D)}{\alpha'(x) + N'(x; D)} \\ &= \frac{\alpha(\cdot, x) + \sum_{n=1}^N \delta(\cdot - Y_n) \delta(x - X_n)}{\alpha'(x) + \sum_{n=1}^N \delta(x - X_n)} \\ &= \left( \frac{\alpha'(x)}{\alpha'(x) + N'(x; D)} \right) \frac{\alpha(\cdot, x)}{\alpha'(x)} + \left( \frac{N'(x; D)}{\alpha'(x) + N'(x; D)} \right) \frac{\bar{N}(\cdot, x; D)}{N'(x; D)}, \end{aligned} \quad (5.27)$$

where  $\alpha'(x) = \int_{\mathcal{Y}} \alpha(y, x) dy$  and  $N'(x; D) = \sum_{n=1}^N \delta(x - X_n)$ .

The conditional distribution when  $\mathcal{X}$  is a Euclidean space has notable differences from its form for a countable set  $\mathcal{X}$ . Specifically, as  $N'(x; D) \in [0, \infty)$ , and in fact will either be zero or tend towards infinity, the coefficients dictating the convex combination of distributions will be zero or one (assuming a non-impulsive model parameter  $\alpha$ ). Thus, the distribution for a given observation  $x$  will be either strictly dependent on either the training data or the prior knowledge regarding  $\theta$ .

#### 5.4.2.3 Training Data PDF, $p_D$

By the Dirichlet process properties,

$$p_{D_{n+1}|D_n, \dots, D_1} = \frac{\alpha(\cdot, \cdot) + \sum_{i=1}^n \delta(\cdot - Y_i) \delta(\cdot - X_i)}{\alpha_0 + N} \quad (5.28)$$

and thus

$$\begin{aligned} p_D(D) &= p(D_1) \prod_{n=2}^N p(D_n|D_{n-1}, \dots, D_1) \\ &= \frac{\alpha(Y_1, X_1)}{\alpha_0} \prod_{n=2}^N \frac{\alpha(Y_n, X_n) + \sum_{i=1}^{n-1} \delta(Y_n - Y_i) \delta(X_n - X_i)}{\alpha_0 + n - 1} \end{aligned} \quad (5.29)$$

It is instructional to find the PDF's for the training output values  $Y$  given the input values  $X$ , as well as the marginal PDF for the input values alone. Observe that the PDF for  $X$  can be represented as

$$\begin{aligned} p_X(X) &= E_\theta [p_{X|\theta}(X|\theta)] = E_\theta \left[ \prod_{n=1}^N \theta'(X_n) \right] \\ &= \frac{\alpha'(X_1)}{\alpha_0} \prod_{n=2}^N \frac{\alpha'(X_n) + \sum_{i=1}^{n-1} \delta(X_n - X_i)}{\alpha_0 + n - 1} \end{aligned} \quad (5.30)$$

since, by the aggregation principle,  $\theta' = \int_Y \theta(y, \cdot) dy$  is a Dirichlet process with parameter function  $\alpha' : \mathcal{X} \mapsto \mathbb{R}^+$ .

Additionally, by the invariance principle, the PDF's for the first-degree marginals are

$$p_{X_n} = \frac{\alpha'}{\alpha_0}. \quad (5.31)$$

which equivalent to  $p_x$ .

PGR: express conditional below using Dir aggregation conditional independence properties?

The conditional distribution of intererest is

$$p_{Y|X}(Y|X) = \frac{\alpha(Y_1, X_1)}{\alpha'(X_1)} \prod_{n=2}^N \frac{\alpha(Y_n, X_n) + \sum_{i=1}^{n-1} \delta(Y_n - Y_i) \delta(X_n - X_i)}{\alpha'(X_n) + \sum_{i=1}^{n-1} \delta(X_n - X_i)} \quad (5.32)$$

Marginalized conditional PDF's for the first and second samples are found. Observe that the marginal distribution for the first  $N - 1$  values of  $Y$  is

$$\begin{aligned} &p_{Y_1, \dots, Y_{N-1}|X}(Y_1, \dots, Y_{N-1}|X) \\ &= \int_Y \frac{\alpha(Y_1, X_1)}{\alpha'(X_1)} \prod_{n=2}^N \frac{\alpha(Y_n, X_n) + \sum_{i=1}^{n-1} \delta(Y_n - Y_i) \delta(X_n - X_i)}{\alpha'(X_n) + \sum_{i=1}^{n-1} \delta(X_n - X_i)} dY_N \\ &= \frac{\alpha(Y_1, X_1)}{\alpha'(X_1)} \prod_{n=2}^{N-1} \frac{\alpha(Y_n, X_n) + \sum_{i=1}^{n-1} \delta(Y_n - Y_i) \delta(X_n - X_i)}{\alpha'(X_n) + \sum_{i=1}^{n-1} \delta(X_n - X_i)} \end{aligned} \quad (5.33)$$

which is independent of  $X_N$ . Repeated integrations and an application of the permutation invariance principle can show that when conditioned on  $X$  any subset of training data values  $Y_1, \dots, Y_N$  will only be dependent on the corresponding values  $X_n$ . The first and second order conditional distributions are

$$P_{Y_n|X_n}(y|x) = \frac{\alpha(y, x)}{\alpha'(x)} = P_{Y|X}(y|x) \quad (5.34)$$

and

$$\begin{aligned} & P_{Y_n, Y_{n'} | X_n, X_{n'}}(y, y' | x, x') \\ &= \frac{\alpha(y, x)\alpha(y', x') + \alpha(y, x)\delta(y - y')\delta(x - x')}{\alpha'(x)\alpha'(x') + \alpha'(x')\delta(x - x')} \end{aligned} \quad (5.35)$$

and the first and second order moments of interest are

$$\mu_{Y_n | X} = \mu_{y | x}(X_n), \quad (5.36)$$

$$E_{Y_n | X} [Y_n^2] = E_{y | x} [y^2](X_n), \quad (5.37)$$

and

$$\begin{aligned} & E_{Y_n | X} [Y_n Y_{n'}] \\ &= \frac{\alpha'(X_n)\mu_{y | x}(X_n)\mu_{y | x}(X_{n'}) + E_{y | x} [y^2](X_n)\delta(X_n - X_{n'})}{\alpha'(X_n) + \delta(X_n - X_{n'})}. \end{aligned} \quad (5.38)$$

PGR: formalize permutation invariance principle???

PGR: Add Y given X equations (with Betas) for discrete case in previous chapters?

PGR: Dirichlet-Multinomial Process perspective

We have the DMP  $\bar{n} \equiv \bar{N}(D)$  with mean and correlation functions

$$\mu_{\bar{n}} = N \frac{\alpha}{\alpha_0} \quad (5.39)$$

and

$$\begin{aligned} & E_{\bar{n}} [\bar{n}(y, x)\bar{n}(y', x')] \\ &= \frac{N}{\alpha_0(\alpha_0 + 1)} [(N - 1)\alpha(y, x)\alpha(y', x') + (\alpha_0 + N)\alpha(y, x)\delta(y - y')\delta(x - x')]. \end{aligned} \quad (5.40)$$

Observe that by the aggregation principle,  $n' = \int_Y \bar{n}(y, \cdot) dy \equiv \sum_{n=1}^N \delta(\cdot - X_n)$  is a DMP over the set  $\mathcal{X}$  with parametrizing function  $\alpha' : \mathcal{X} \mapsto \mathbb{R}^+$ .

Additionally, the 1-dimensional subsets conditioned on the marginalized DMP are characterized as

$$\frac{\bar{n}(\cdot, x)}{\delta(0)} \Big| n'(x) \sim \text{DMP} \left( \frac{n'(x)}{\delta(0)}, \frac{\alpha(\cdot, x)}{\delta(0)} \right) \quad (5.41)$$

PGR: add proof???

## 5.5 Applications: General Model

PGR: COPIED, incomplete

$$\begin{aligned}
 E_{y|x,D} [\mathcal{L}(h,y)] &= \int_{\mathcal{Y}} \mathcal{L}(h,y) p_{y|x,D}(y|x,D) dy \\
 &= \frac{\int_{\mathcal{Y}} \alpha(y,x) \mathcal{L}(h,y) dy + \int_{\mathcal{Y}} \bar{N}(y,x;D) \mathcal{L}(h,y) dy}{\alpha'(x) + N'(x;D)} \\
 &= \frac{\int_{\mathcal{Y}} \alpha(y,x) \mathcal{L}(h,y) dy + \sum_{n=1}^N \delta(x - X_n) \mathcal{L}(h, D_n)}{\alpha'(x) + N'(x;D)} \\
 &= \left( \frac{\alpha'(x)}{\alpha'(x) + N'(x;D)} \right) E_{y|x} [\mathcal{L}(h,y)] + \left( \frac{N'(x;D)}{\alpha'(x) + N'(x;D)} \right) \frac{\sum_{n=1}^N \delta(x - X_n) \mathcal{L}(h, Y_n)}{\sum_{n=1}^N \delta(x - X_n)}. 
 \end{aligned} \tag{5.42}$$

### 5.5.1 Regression: the Squared-Error Loss

$$\mathcal{L}(h,y) = (h - y)^2. \tag{5.43}$$

Now we choose for the regression function to map to  $\mathcal{H} = \mathcal{Y} = \mathbb{R}$ .

$$\begin{aligned}
 \mathcal{R}(f) &= E_{\theta} \left[ E_{D|\theta} \left[ E_{y|x|\theta} \left[ (f(x;D) - y)^2 \right] \right] \right] \\
 &= E_{x,\theta} \left[ E_{y|x,\theta} [(y - \mu_{y|x,\theta})^2] \right] + E_{\theta} \left[ E_{x,D|\theta} \left[ (f(x;D) - \mu_{y|x,\theta})^2 \right] \right] \\
 &= E_{x,\theta} [\Sigma_{y|x,\theta}] + E_{\theta} \left[ E_{x,D|\theta} \left[ (f(x;D) - \mu_{y|x,\theta})^2 \right] \right]
 \end{aligned} \tag{5.44}$$

### 5.5.1.1 Optimal Learner

The optimal function is the expected value of the output conditional PDF,

$$\begin{aligned}
 f^*(x; D) &= \mu_{y|x,D} = E_{\theta|x,D} [\mu_{y|x,\theta}] \\
 &= \frac{\int_{\mathcal{Y}} y(\alpha(y,x) + N(y,x; D)) dy}{\alpha'(x) + N'(x; D)} \\
 &= \left( \frac{\alpha'(x)}{\alpha'(x) + N'(x; D)} \right) \int_{\mathcal{Y}} y \frac{\alpha(y,x)}{\alpha'(x)} dy \\
 &\quad + \left( \frac{N'(x; D)}{\alpha'(x) + N'(x; D)} \right) \frac{\sum_{n=1}^N \delta(x - X_n) Y_n}{\sum_{n=1}^N \delta(x - X_n)} \\
 &= \left( \frac{\alpha'(x)}{\alpha'(x) + N'(x; D)} \right) \mu_{y|x} \\
 &\quad + \left( \frac{N'(x; D)}{\alpha'(x) + N'(x; D)} \right) \frac{\sum_{n=1}^N \delta(x - X_n) Y_n}{\sum_{n=1}^N \delta(x - X_n)}.
 \end{aligned} \tag{5.45}$$

### 5.5.1.2 Minimum Risk

Generalizing from the basic model discussion, we again have

$$\begin{aligned}
 \mathcal{R}^* &= E_{x,D} [\Sigma_{y|x,D}] = E_{x,\bar{n}} [\Sigma_{y|x,\bar{n}}] \\
 &= E_{x,\theta} [\Sigma_{y|x,\theta}] + E_{x,D} [C_{\theta|x,D} [\mu_{y|x,\theta}]],
 \end{aligned} \tag{5.46}$$

where we choose to perform the expectation over  $\bar{n}$ .

The conditional variance is now

$$\Sigma_{y|x,\bar{n}} = E_{y|x,\bar{n}} [y^2] - \mu_{y|x,\bar{n}}^2. \tag{5.47}$$

and the two terms are independently evaluated.

$$\begin{aligned}
 E_{x,D} [E_{y|x,D} [y^2]] &= E_{x,\bar{n}} [E_{y|x,\bar{n}} [y^2]] \\
 &= E_y [y^2] = \int_{\mathcal{Y}} y^2 \int_{\mathcal{X}} \frac{\alpha(y,x)}{\alpha_0} dx dy \\
 &= E_x [E_{y|x} [y^2]] = \int_{\mathcal{X}} \frac{\alpha'(x)}{\alpha_0} \int_{\mathcal{Y}} y^2 \frac{\alpha(y,x)}{\alpha'(x)} dy dx.
 \end{aligned} \tag{5.48}$$

$$\begin{aligned}
 E_{x,D} [\mu_y^2 |_{x,D}] &= E_x E_{D|x} [\mu_y^2 |_{x,D}] \\
 &= E_x \left[ E_{D|x} \left[ \left( \frac{\alpha'(x)\mu_{y|x} + \sum_{n=1}^N Y_n \delta(x - X_n)}{\alpha'(x) + \sum_{n=1}^N \delta(x - X_n)} \right)^2 \right] \right] \\
 &= E_x \left[ E_D \left[ \frac{\alpha_0 \left( \alpha'(x)\mu_{y|x} + \sum_{n=1}^N Y_n \delta(x - X_n) \right)^2}{\alpha'(x) \left( \alpha'(x) + \sum_{n=1}^N \delta(x - X_n) \right) (\alpha_0 + N)} \right] \right] \\
 &= E_x \left[ E_X \left[ \frac{\alpha_0 E_{Y|x} \left[ \left( \alpha'(x)\mu_{y|x} + \sum_{n=1}^N Y_n \delta(x - X_n) \right)^2 \right]}{\alpha'(x) \left( \alpha'(x) + \sum_{n=1}^N \delta(x - X_n) \right) (\alpha_0 + N)} \right] \right]
 \end{aligned} \tag{5.49}$$

Evaluating the expectation over Y given X, we have

$$\begin{aligned}
 E_{Y|X} \left[ \left( \alpha'(x)\mu_{y|x} + \sum_{n=1}^N Y_n \delta(x - X_n) \right)^2 \right] &= \alpha'(x)^2 \mu_{y|x}^2 + 2\alpha'(x)\mu_{y|x} \sum_{n=1}^N \mu_{y|x}(X_n) \delta(x - X_n) \\
 &\quad + \sum_{n=1}^N E_{y|x} [y^2](X_n) \delta(x - X_n)^2 \\
 &\quad + \sum_{n \neq n'} \frac{\alpha'(X_n)\mu_{y|x}(X_n)\alpha'(X_{n'})\mu_{y|x}(X_{n'}) + \alpha'(X_n)E_{y|x}[y^2](X_n)\delta(X_n - X_{n'})}{\alpha'(X_n)\alpha'(X_{n'}) + \alpha'(X_n)\delta(X_n - X_{n'})} \\
 &\quad \delta(x - X_n)\delta(x - X_{n'}) \\
 &= \dots \\
 &= \alpha'(x)^2 \mu_{y|x}^2 + 2\alpha'(x)\mu_{y|x}^2 \sum_{n=1}^N \delta(x - X_n) + E_{y|x} [y^2] \sum_{n=1}^N \delta(x - X_n)^2 \\
 &\quad + \frac{\alpha'(x)\mu_{y|x}^2 + E_{y|x}[y^2]\delta(0)}{\alpha'(x) + \delta(0)} \sum_{n \neq n'} \delta(x - X_n)\delta(x - X_{n'}) \\
 &= \dots \\
 &= \frac{\alpha'(x) + \sum_{n=1}^N \delta(x - X_n)}{\alpha'(x) + \delta(0)} \\
 &\quad \left( E_{y|x} [y^2]\delta(0) \sum_{n=1}^N \delta(x - X_n) + \alpha'(x)\mu_{y|x}^2 \left( \alpha'(x) + \delta(0) + \sum_{n=1}^N \delta(x - X_n) \right) \right)
 \end{aligned} \tag{5.50}$$

PGR: Y given X PDFs, moments??? In PDF section, or in Appendix?

Plugging,

$$\begin{aligned}
 & E_{x,D} \left[ \mu_{y|x,D}^2 \right] \tag{5.51} \\
 &= E_x \left[ E_X \left[ \frac{\alpha_0 E_{Y|X} \left[ (\alpha'(x) \mu_{y|x} + \sum_{n=1}^N Y_n \delta(x - X_n))^2 \right]}{\alpha'(x) (\alpha'(x) + \sum_{n=1}^N \delta(x - X_n)) (\alpha_0 + N)} \right] \right] \\
 &= E_x \left[ \frac{\alpha_0 E_X \left[ E_{y|x} [y^2] \delta(0) \sum_{n=1}^N \delta(x - X_n) + \alpha'(x) \mu_{y|x}^2 (\alpha'(x) + \delta(0) + \sum_{n=1}^N \delta(x - X_n)) \right]}{\alpha'(x) (\alpha'(x) + \delta(0)) (\alpha_0 + N)} \right]
 \end{aligned}$$

Evaluating the expectation over  $X$ ,

$$\begin{aligned}
 & E_X \left[ E_{y|x} [y^2] \delta(0) \sum_{n=1}^N \delta(x - X_n) + \alpha'(x) \mu_{y|x}^2 \left( \alpha'(x) + \delta(0) + \sum_{n=1}^N \delta(x - X_n) \right) \right] \\
 &= E_{y|x} [y^2] \delta(0) N \frac{\alpha'(x)}{\alpha_0} + \alpha'(x) \mu_{y|x}^2 \left( \alpha'(x) + \delta(0) + N \frac{\alpha'(x)}{\alpha_0} \right) \\
 &= \frac{\alpha'(x)}{\alpha_0} \left( E_{y|x} [y^2] \delta(0) N + \mu_{y|x}^2 (\alpha_0 \alpha'(x) + \alpha_0 \delta(0) + N \alpha'(x)) \right)
 \end{aligned}$$

Plugging,

$$\begin{aligned}
 & E_{x,D} \left[ \mu_{y|x,D}^2 \right] \tag{5.52} \\
 &= E_x \left[ \frac{E_{y|x} [y^2] \delta(0) N + \mu_{y|x}^2 (\alpha_0 \alpha'(x) + \alpha_0 \delta(0) + N \alpha'(x))}{(\alpha'(x) + \delta(0)) (\alpha_0 + N)} \right]
 \end{aligned}$$

Combining with the second moment produces the risk,

$$\begin{aligned}
 \mathcal{R}^* &= E_{x,D} \left[ E_{y|x,D} [y^2] - \mu_{y|x,D}^2 \right] \tag{5.53} \\
 &= E_x \left[ \frac{\alpha_0 \alpha'(x) + \alpha_0 \delta(0) + N \alpha'(x)}{(\alpha_0 + N) (\alpha'(x) + \delta(0))} \Sigma_{y|x} \right] \\
 &= E_x \left[ \frac{P_x(x) + (\alpha_0 + N)^{-1} \delta(0)}{P_x(x) + \alpha_0^{-1} \delta(0)} \Sigma_{y|x} \right]
 \end{aligned}$$

PGR: Discuss Dirac deltas!!!

PGR: DMP PERSPECTIVE???

To perform the expectation over the Dirichlet-Multinomial process  $\bar{n} \sim \text{DMP}(\alpha)$ , split the expectation into an expectation over the marginal DMP  $n' \sim \text{DMP}(\alpha')$  and a conditional expectation over  $\bar{n}$  given  $n'$ . The characterization of the conditional DMP is found in Appendix A.3.

$$\begin{aligned}
 & E_{x,\bar{n}} \left[ \mu_{y|x,\bar{n}}^2 \right] \tag{5.54} \\
 &= E_x \left[ E_{\bar{n}|x} \left[ \left( \frac{\alpha'(x)\mu_{y|x} + \int_{\mathcal{Y}} y\bar{n}(y,x)dy}{\alpha'(x) + n'(x)} \right)^2 \right] \right] \\
 &= E_x E_{n'} \left[ \frac{\alpha_0 E_{\bar{n}|n'} \left[ (\alpha'(x)\mu_{y|x} + \int_{\mathcal{Y}} y\bar{n}(y,x)dy)^2 \right]}{\alpha'(x)(\alpha'(x) + n'(x))(\alpha_0 + N)} \right]
 \end{aligned}$$

Evaluating the conditional expectation,

$$\begin{aligned}
 & E_{\bar{n}|n'} \left[ \left( \alpha'(x)\mu_{y|x} + \int_{\mathcal{Y}} y\bar{n}(y,x)dy \right)^2 \right] \tag{5.55} \\
 &= \alpha'(x)^2 \mu_{y|x}^2 + 2\alpha'(x)\mu_{y|x} \int_{\mathcal{Y}} y\delta(0) \frac{n'(x)}{\delta(0)} \frac{\delta(0)^{-1}\alpha(y,x)}{\delta(0)^{-1}\alpha'(x)} dy \\
 &\quad + \delta(0)^2 \frac{\delta(0)^{-1}n'(x)}{(\delta(0)^{-1}\alpha'(x))(\delta(0)^{-1}\alpha'(x) + 1)} \int_{\mathcal{Y}} \int_{\mathcal{Y}} yy' \left[ \left( \frac{n'(x)}{\delta(0)} - 1 \right) \frac{\alpha(y,x)}{\delta(0)} \frac{\alpha(y',x)}{\delta(0)} \right. \\
 &\quad \left. + \left( \frac{\alpha'(x)}{\delta(0)} + \frac{n'(x)}{\delta(0)} \right) \frac{\alpha(y,x)}{\delta(0)} \delta(y-y') \right] dy dy' \\
 &= \alpha'(x)^2 \mu_{y|x}^2 + 2\alpha'(x)n'(x)\mu_{y|x}^2 \\
 &\quad + \frac{n'(x)}{\alpha'(x)(\alpha'(x) + \delta(0))} \left[ (n'(x) - \delta(0))\alpha'(x)^2 \mu_{y|x}^2 + \delta(0)(\alpha'(x) + n'(x))\alpha'(x) E_{y|x} [y^2] \right] \\
 &= \frac{\alpha'(x) + n'(x)}{\alpha'(x) + \delta(0)} \left[ \mu_{y|x}^2 \alpha'(x)(\alpha'(x) + n'(x) + \delta(0)) + E_{y|x} [y^2] \delta(0) n'(x) \right]
 \end{aligned}$$

Plugging,

$$\begin{aligned}
 & E_{x,\bar{n}} \left[ \mu_{y|x,\bar{n}}^2 \right] \tag{5.56} \\
 &= E_x \left[ \frac{\alpha_0 E_{n'} \left[ \mu_{y|x}^2 \alpha'(x)(\alpha'(x) + n'(x) + \delta(0)) + E_{y|x} [y^2] \delta(0) n'(x) \right]}{\alpha'(x)(\alpha'(x) + \delta(0))(\alpha_0 + N)} \right] \\
 &= E_x \left[ \frac{\mu_{y|x}^2 (\alpha_0 \alpha'(x) + N \alpha'(x) + \delta(0) \alpha_0) + E_{y|x} [y^2] \delta(0) N}{(\alpha_0 + N)(\alpha'(x) + \delta(0))} \right]
 \end{aligned}$$

Combining produces the risk,

$$\begin{aligned}
 \mathcal{R}^* &= E_x \left[ \frac{\alpha_0 \alpha'(x) + \alpha_0 \delta(0) + N \alpha'(x)}{(\alpha_0 + N)(\alpha'(x) + \delta(0))} \Sigma_{y|x} \right] \tag{5.57} \\
 &= E_x \left[ \frac{P(x) + (\alpha_0 + N)^{-1} \delta(0)}{P(x) + \alpha_0^{-1} \delta(0)} \Sigma_{y|x} \right]
 \end{aligned}$$

# Appendix A

PGR: notation/formatting

PGR: remove beta set arguments by introducing alpha sub z?

## A.1 Dirichlet random process conditioned on its aggregation

This section details an important property of Dirichlet distributed random processes. The following development first considers Dirichlet random processes over a countable domain and then generalizes for continuous-domain Dirichlet processes.

First, define the PDF of a Dirichlet aggregation [7]. Let the random process  $\theta \in \Theta = \mathcal{P}(\mathcal{Y})$  be Dirichlet over the countable set  $\mathcal{Y}$  with parameterizing function  $\alpha \in \mathbb{R}^{+\mathcal{Y}}$ . Define an arbitrary partition of  $\mathcal{Y}$ :  $\{\dots, \mathcal{S}(z), \dots\}$ ,  $z \in \mathcal{Z}$ ; the aggregation  $\theta' \in \mathcal{P}(\mathcal{Z})$ ,  $\theta'(z) \equiv \sum_{y \in \mathcal{S}(z)} \theta(y)$  is thus also Dirichlet and has a parameterizing function  $\alpha' \in \mathbb{R}^{+\mathcal{Z}}$ ,  $\alpha'(z) \equiv \sum_{y \in \mathcal{S}(z)} \alpha(y)$ .

The PDF of the original random process  $\theta$  conditioned on its aggregation  $\theta'$  can be formulated as

$$\begin{aligned} p_{\theta|\theta'}(\theta|\theta') &= \frac{\beta(\alpha') \prod_{y \in \mathcal{Y}} \theta(y)^{\alpha(y)-1}}{\beta(\alpha) \prod_{z \in \mathcal{Z}} \theta'(z)^{\alpha'(z)-1}} & (A.1) \\ &= \prod_{z \in \mathcal{Z}} \left[ \beta\left(\{\alpha(y) : y \in \mathcal{S}(z)\}\right)^{-1} \frac{\prod_{y \in \mathcal{S}(z)} \theta(y)^{\alpha(y)-1}}{\theta'(z)^{\alpha'(z)-1}} \right] \\ &= \prod_{z \in \mathcal{Z}} \left[ \frac{\theta'(z)^{1-|\mathcal{S}(z)|}}{\beta\left(\{\alpha(y) : y \in \mathcal{S}(z)\}\right)} \prod_{y \in \mathcal{S}(z)} \left( \frac{\theta(y)}{\theta'(z)} \right)^{\alpha(y)-1} \right], \end{aligned}$$

which is defined for  $\left\{ \theta \in \mathbb{R}_{\geq 0}^{\mathcal{Y}} : \sum_{y \in \mathcal{S}(z)} \theta(y) = \theta'(z), \quad \forall z \in \mathcal{Z} \right\}$ .

Observe that the partitioned segments are conditionally independent; introduce subscript notation to refer to the function segment  $\theta_z = \{\theta(y) : y \in \mathcal{S}(z)\}$ . The PDF  $p_{\theta|\theta'}$  can now be decomposed as  $p_{\theta|\theta'}(\dots, \theta_z, \dots | \theta') = \prod_{z \in \mathcal{Z}} p_{\theta_z|\theta'(z)}(\theta_z | \theta'(z))$

Next, normalize the segments of  $\theta$  to form  $\tilde{\theta} = (\dots, \tilde{\theta}_z, \dots)$ , where  $\tilde{\theta}_z \equiv \theta_z / \theta'(z)$ , and formulate the conditional PDF

$$\begin{aligned} p_{\tilde{\theta}|\theta'}(\tilde{\theta} | \theta') &= \prod_{z \in \mathcal{Z}} \left[ \frac{\prod_{y \in \mathcal{S}(z)} \tilde{\theta}_z(y)^{\alpha(y)-1}}{\beta(\{\alpha(y) : y \in \mathcal{S}(z)\})} \right] \\ &= \prod_{z \in \mathcal{Z}} \text{Dir}(\tilde{\theta}_z; \{\alpha(y) : y \in \mathcal{S}(z)\}). \end{aligned} \quad (\text{A.2})$$

which is defined for  $\tilde{\theta} \in \prod_{z \in \mathcal{Z}} \{\tilde{\theta}_z \in \mathcal{P}(\mathcal{S}(z))\}$ . Thus after conditioning, the normalized segments  $\tilde{\theta}_z$  are Dirichlet distributed, independent of one another, and independent of the aggregation  $\theta'$ .

PGR: discuss transform Jacobian and dimensionality?

This principle holds for continuous-domain Dirichlet processes  $\theta$  as well - the segments  $\tilde{\theta}_z$  are now continuous-domain Dirichlet processes.

## A.2 Multinomial Distribution Properties

### A.2.1 Aggregation

A characteristic of a Multinomial random process is that its aggregations are also Multinomial [9]. Consider a random process  $\bar{n} \sim \text{Multi}(N, \theta)$  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 distributed as  $n' \sim \text{Multi}(N, \theta')$  with parameterizing function  $\theta'(z) = \sum_{y \in \mathcal{S}(z)} \theta(y)$ .

To prove this principle, define the subset  $\tilde{\mathcal{N}} = \{\bar{n} \in \mathcal{N} : \sum_{y \in \mathcal{S}(z)} \bar{n}(y) = n'(z), \quad \forall z \in$

$\mathcal{Z} \} \subseteq \bar{\mathcal{N}}$ , where the original random process  $\bar{n} \in \bar{\mathcal{N}}$ . Next, observe

$$\begin{aligned} P_{n'}(n') &= \sum_{\bar{n} \in \bar{\mathcal{N}}} P_{\bar{n}}(\bar{n}) \\ &= \mathcal{M}(n') \prod_{z \in \mathcal{Z}} \sum_{\substack{n'(z) = \\ \sum_{y \in \mathcal{S}(z)} \bar{n}(y)}} \mathcal{M}(\{\bar{n}(y) : y \in \mathcal{S}(z)\}) \prod_{y \in \mathcal{S}(z)} \theta(y)^{\bar{n}(y)} \\ &= \mathcal{M}(n') \prod_{z \in \mathcal{Z}} \theta'(z)^{n'(z)} = \text{Multi}(n'; N, \theta') , \end{aligned} \tag{A.3}$$

where the multinomial theorem [8] has been used.

### A.2.2 Conditioned on its Aggregation

If the multinomial random process  $\bar{n}$  is conditioned on its aggregation over the partition  $\{\dots, \mathcal{S}(z), \dots\}$ ,  $z \in \mathcal{Z}$ , the distinct segments  $\bar{n}(y)$ ,  $y \in \mathcal{S}(z)$  become independent multinomial random processes,

$$\begin{aligned} P_{\bar{n}|n'}(\bar{n}|n') &= \frac{\mathcal{M}(\bar{n}) \prod_{y \in \mathcal{Y}} \theta(y)^{\bar{n}(y)}}{\mathcal{M}(n') \prod_{z \in \mathcal{Z}} \theta'(z)^{n'(z)}} \\ &= \prod_{z \in \mathcal{Z}} \left[ \mathcal{M}(\{\bar{n}(y) : y \in \mathcal{S}(z)\}) \prod_{y \in \mathcal{S}(z)} \left( \frac{\theta(y)}{\theta'(z)} \right)^{\bar{n}(y)} \right] \\ &= \prod_{z \in \mathcal{Z}} \left[ \mathcal{M}(\{\bar{n}(y) : y \in \mathcal{S}(z)\}) \prod_{y \in \mathcal{S}(z)} \tilde{\theta}_z(y)^{\bar{n}(y)} \right] \\ &= \prod_{z \in \mathcal{Z}} \text{Multi} \left( \{\bar{n}(y) : y \in \mathcal{S}(z)\}; n'(z), \{\tilde{\theta}_z(y) : y \in \mathcal{S}(z)\} \right) , \end{aligned} \tag{A.4}$$

on the domain  $\{\bar{n} \in \mathbb{Z}_{\geq 0}^{\mathcal{Y}} : \sum_{y \in \mathcal{S}(z)} \bar{n}(y) = n'(z), \quad \forall z \in \mathcal{Z}\}$ . Observe that the segment over the set  $\mathcal{S}(z)$  sums to  $n'(z)$  and is parameterized by the normalized segment  $\tilde{\theta}_z \equiv \{\theta(y)/\theta'(z) : y \in \mathcal{S}(z)\}$ .

## A.3 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 [9]. 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  $\alpha'(z) = \sum_{y \in \mathcal{S}(z)} \alpha(y)$ .

It can be shown that 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

$$\begin{aligned} P_{\bar{n}|n'}(\bar{n}|n') &= \frac{\mathcal{M}(\bar{n})\beta(\alpha)^{-1}\beta(\alpha + \bar{n})}{\mathcal{M}(n')\beta(\alpha')^{-1}\beta(\alpha' + n')} \\ &= \left( \prod_{z \in \mathcal{Z}} \frac{\Gamma(\alpha'(z) + n'(z))}{n'(z)!\Gamma(\alpha'(z))} \right)^{-1} \left( \prod_{y \in \mathcal{Y}} \frac{\Gamma(\alpha(y) + \bar{n}(y))}{\bar{n}(y)!\Gamma(\alpha(y))} \right) \\ &= \prod_{z \in \mathcal{Z}} \left[ \frac{n'(z)!\Gamma(\alpha'(z))}{\Gamma(\alpha'(z) + n'(z))} \prod_{y \in \mathcal{S}(z)} \frac{\Gamma(\alpha(y) + \bar{n}(y))}{\bar{n}(y)!\Gamma(\alpha(y))} \right] \\ &= \prod_{z \in \mathcal{Z}} \text{DM} \left( \{\bar{n}(y) : y \in \mathcal{S}(z)\}; n'(z), \{\alpha(y) : y \in \mathcal{S}(z)\} \right), \end{aligned} \quad (\text{A.5})$$

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

## A.4 First and Second moments of a Dirichlet Process

In this section, it is shown that the expected value of a Dirichlet process  $\theta \sim \text{DP}(\alpha)$  is

$$\mu_\theta = \frac{\alpha}{\alpha_0}, \quad (\text{A.6})$$

where  $\alpha_0 = \int_{\mathcal{Y}} \alpha(y) dy$ .

The defining characteristic of Dirichlet processes is that their aggregations are also Dirichlet. Define the partition of  $\mathcal{Y}$ ,  $\{\mathcal{S}(y), \mathcal{S}^c(y)\}$  where  $\mathcal{S}(y) = \{t \in \mathcal{Y} : t \leq y\}$ . The transform random variable  $\theta'(y) \equiv \int_{-\infty}^y \theta(t) dt$  is thus a Beta random variable with parameters  $\lambda = \int_{-\infty}^y \alpha(t) dt$  and  $\lambda^c = \int_y^\infty \alpha(t) dt$ . Using the formula for the expected value of a beta random variable [13], note that

$$\begin{aligned} \mu_{\theta'} &= \frac{\lambda}{\lambda + \lambda^c} \\ &= \frac{\int_{-\infty}^y \alpha(t) dt}{\alpha_0} = \int_{-\infty}^y \mu_\theta(t) dt. \end{aligned} \quad (\text{A.7})$$

Differentiating with respect to  $y$ , we have the expected value of the DP.

Next, the correlation function is shown to be

$$E_{\theta} [\theta(y_1)\theta(y_2)] = \frac{\alpha(y_1)\alpha(y_2) + \alpha(y_1)\delta(y_1 - y_2)}{\alpha_0(\alpha_0 + 1)}. \quad (\text{A.8})$$

First, assume  $y_2 \geq y_1$  and define a new partition of  $\mathcal{Y}$ ,  $\{(-\infty, y_1], (y_1, y_2], (y_2, \infty)\}$ . By the aggregation property, the random triplet  $(\int_{-\infty}^{y_1} \theta(t)dt, \int_{y_1}^{y_2} \theta(t)dt, \int_{y_2}^{\infty} \theta(t)dt)$  is Dirichlet with parameters  $(\int_{-\infty}^{y_1} \alpha(t)dt, \int_{y_1}^{y_2} \alpha(t)dt, \int_{y_2}^{\infty} \alpha(t)dt)$ .

Define the function

$$\begin{aligned} g(t_1, t_2) &= E_{\theta} \left[ \int_{-\infty}^{y_1} \theta(t_1)dt_1 \int_{-\infty}^{y_2} \theta(t_2)dt_2 \right] \\ &= E_{\theta} \left[ \left( \int_{-\infty}^{y_1} \theta(t_1)dt_1 \right)^2 + \left( \int_{-\infty}^{y_1} \theta(t_1)dt_1 \right) \left( \int_{y_1}^{y_2} \theta(t_2)dt_2 \right) \right] \\ &= \frac{\left( \int_{-\infty}^{y_1} \alpha(t_1)dt_1 \right) \left( 1 + \int_{-\infty}^{y_1} \alpha(t_1)dt_1 \right) + \left( \int_{-\infty}^{y_1} \alpha(t_1)dt_1 \right) \left( \int_{y_1}^{y_2} \alpha(t_2)dt_2 \right)}{\alpha_0(\alpha_0 + 1)} \\ &= \frac{\left( \int_{-\infty}^{y_1} \alpha(t_1)dt_1 \right) \left( \int_{-\infty}^{y_2} \alpha(t_2)dt_2 \right) + \left( \int_{-\infty}^{y_1} \alpha(t_1)dt_1 \right)}{\alpha_0(\alpha_0 + 1)} \quad \forall y_2 \geq y_1. \end{aligned} \quad (\text{A.9})$$

Following the same steps provides the values of  $g$  for  $t_2 \leq t_1$ ; the combined formula can be given as

$$g(t_1, t_2) = \frac{\left( \int_{-\infty}^{y_1} \alpha(t_1)dt_1 \right) \left( \int_{-\infty}^{y_2} \alpha(t_2)dt_2 \right) + \left( \int_{-\infty}^{\min(y_1, y_2)} \alpha(t_1)dt_1 \right)}{\alpha_0(\alpha_0 + 1)}. \quad (\text{A.10})$$

Finally,

$$\begin{aligned} E_{\theta} [\theta(y_1)\theta(y_2)] &= \frac{d^2}{dt_1 dt_2} g(t_1, t_2) \\ &= \frac{\frac{d}{dt_2} \left[ \alpha(y_1) \left( \int_{-\infty}^{y_2} \alpha(t_2)dt_2 \right) + u(t_2 - t_1)\alpha(\min(t_1, t_2)) \right]}{\alpha_0(\alpha_0 + 1)} \\ &= \frac{\alpha(y_1)\alpha(y_2) + \alpha(y_1)\delta(y_1 - y_2)}{\alpha_0(\alpha_0 + 1)}. \end{aligned} \quad (\text{A.11})$$

## A.5 Proof: Continuous Model Posterior Distribution is Dirichlet Process

In this section, it is shown that if the model  $\theta \sim \text{DP}(\alpha)$  is a Dirichlet process, then the model conditioned on the training data  $D$  is also a Dirichlet process with parameterizing function

$$\alpha + \sum_{n=1}^N \delta(\cdot - D_n).$$

The defining characteristic of Dirichlet processes is that their aggregations are also Dirichlet. Consider a DP 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  $\theta'(z) \equiv \int_{\mathcal{S}(z)} \theta(y) dy$  is necessarily Dirichlet with parameterizing function  $\alpha'(z) \equiv \int_{\mathcal{S}(z)} \alpha(y) dy$ .

To prove the hypothesis, it is required that

$$\theta' | D \sim \text{Dir}(\alpha' + \bar{N}(D)), \quad (\text{A.12})$$

where  $\bar{N}(z; D) = \int_{\mathcal{S}(z)} \sum_{n=1}^N \delta(y - D_n) dy = \sum_{n=1}^N \chi(D_n; \mathcal{S}(z))$ .  $\chi$  is the indicator function

$$\chi(x; S) = \begin{cases} 1 & \text{if } x \in S, \\ 0 & \text{if } x \notin S. \end{cases} \quad (\text{A.13})$$

### PGR: SUFFICIENT STATISTIC FOR D!!!?

To prove the hypothesis, exploit the results of Appendix A.1 to represent the training data distribution conditioned on the aggregation  $\theta'$ . The conditional distribution of interest is

$$\begin{aligned} p_{D|\theta'}(D|\theta') &= E_{\theta|\theta'} [p_{D|\theta}(D|\theta)] = E_{\theta|\theta'} \left[ \prod_{n=1}^N \theta(D_n) \right] \\ &= \prod_{z \in \mathcal{Z}} E_{\theta_z|\theta'(z)} \left[ \prod_{n=1}^N \theta_z(D_n)^{\chi(D_n; \mathcal{S}(z))} \right] \\ &= \left( \prod_{z \in \mathcal{Z}} \prod_{n=1}^N \theta'(z)^{\chi(D_n; \mathcal{S}(z))} \right) \prod_{z \in \mathcal{Z}} E_{\tilde{\theta}_z} \left[ \prod_{n=1}^N \tilde{\theta}_z(D_n)^{\chi(D_n; \mathcal{S}(z))} \right] \\ &= \left( \prod_{z \in \mathcal{Z}} \theta'(z)^{\bar{N}(z; D)} \right) \prod_{z \in \mathcal{Z}} E_{\tilde{\theta}_z} \left[ \prod_{n=1}^N \tilde{\theta}_z(D_n)^{\chi(D_n; \mathcal{S}(z))} \right], \end{aligned} \quad (\text{A.14})$$

Observe that the dependency on  $\theta'$  is polynomial. The training data marginal distribution is

$$\begin{aligned} p_D(D) &= E_{\theta'} \left[ \prod_{z \in \mathcal{Z}} \theta'(z)^{\bar{N}(z; D)} \right] \prod_{z \in \mathcal{Z}} E_{\tilde{\theta}_z} \left[ \prod_{n=1}^N \tilde{\theta}_z(D_n)^{\chi(D_n; \mathcal{S}(z))} \right] \\ &= \frac{\beta(\alpha' + \bar{N}(D))}{\beta(\alpha')} \prod_{z \in \mathcal{Z}} E_{\tilde{\theta}_z} \left[ \prod_{n=1}^N \tilde{\theta}_z(D_n)^{\chi(D_n; \mathcal{S}(z))} \right] \end{aligned} \quad (\text{A.15})$$

and thus the distribution of interest is

$$\begin{aligned} p_{\theta'|D}(\theta'|D) &= \frac{\prod_{z \in \mathcal{Z}} \theta'(z)^{\alpha'(z) + \bar{N}(z; D) - 1}}{\beta(\alpha' + \bar{N}(D))} \\ &= \text{Dir}(\theta'; \alpha' + \bar{N}(\cdot; D)). \end{aligned} \quad (\text{A.16})$$

This proves the hypothesis.

## A.6 The Dirichlet-Multinomial Process

PGR: multinomial process too???

### A.6.1 Definition

This section introduces a new random process, referred to as the Dirichlet-Multinomial process (DMP). It is the generalization of the Dirichlet-Multinomial distribution for i.i.d. samples drawn from a PDF; the underlying distribution is characterized by a Dirichlet process with parameter  $\alpha$ . The Dirichlet-Multinomial process assumes functions from the set  $\{\bar{n} \in \mathbb{R}_{\geq 0}^{\mathcal{Y}} : \int_{\mathcal{Y}} \bar{n}(y) dy = N\}$  and is parameterized by a function  $\alpha : \mathcal{Y} \mapsto \mathbb{R}^+$ .

Analogous to the Dirichlet and Dirichlet-Multinomial distributions for countable spaces, the Dirichlet-Multinomial process inherits the aggregation property from the Dirichlet process prior. That is, for a Dirichlet-Multinomial process  $\bar{n} \in \{\bar{n} \in \mathbb{R}_{\geq 0}^{\mathcal{Y}} : \int_{\mathcal{Y}} \bar{n}(y) dy = N\}$  and a partition of  $\mathcal{Y}$ ,  $\{\dots, \mathcal{S}(z), \dots\}$ ,  $z \in \mathcal{Z}$ , the transformed random process  $n'(z) \equiv \int_{\mathcal{S}(z)} \bar{n}(y) dy$  is necessarily Dirichlet-Multinomial with parameterizing function  $\alpha'(z) \equiv \int_{\mathcal{S}(z)} \alpha(y) dy$ .

PGR: tilde not prime?

### A.6.2 Proof that $\sum_{n=1}^N \delta(y - D_n)$ is a DMP

Next, it is demonstrated that the random process  $\bar{n}(y) \equiv \bar{N}(y; D) = \sum_{n=1}^N \delta(y - D_n)$  is a DMP, given that  $p_{D|\theta}(D|\theta) = \prod_{n=1}^N \theta(D_n)$  and  $\theta \sim DP(\alpha)$ .

Observe that  $n'(z) \equiv \sum_{n=1}^N \chi(D_n; \mathcal{S}(z))$ , where  $\chi$  is the indicator function

$$\chi(x; S) = \begin{cases} 1 & \text{if } x \in S, \\ 0 & \text{if } x \notin S. \end{cases} \quad (\text{A.17})$$

and note that  $P(\chi(D_n; \mathcal{S}(z)) = 1 | \theta) = \int_{\mathcal{S}(z)} \theta(y) dy$ . As such,  $n'$  conditioned on the model  $\theta$  is characterized by a multinomial distribution

$$P_{n'|\theta}(n'|\theta) = \mathcal{M}(n') \prod_{z \in \mathcal{Z}} \left( \int_{\mathcal{S}(z)} \theta(y) dy \right)^{n'(z)} = \text{Multi}(n'; N, \theta'(z)) , \quad (\text{A.18})$$

where  $\theta'(z) \equiv \int_{\mathcal{S}(z)} \theta(y) dy, z \in \mathcal{Z}$ .

By the aggregation property of the Dirichlet process  $\theta$ , the parameters of this multinomial distribution are characterized as  $\theta' \sim \text{Dir}(\alpha')$ , and thus  $\bar{n}$  is drawn from a Dirichlet-Multinomial PMF with the same parameters  $\alpha'$ . As this holds for any partition of  $\mathcal{Y}$ ,  $\bar{n}$  is a Dirichlet-Multinomial Process.

### A.6.3 Mean and Correlation Functions

In this subsection the mean and correlation functions of a DMP are expressed. The mean function is

$$\begin{aligned} \mu_{\bar{n}}(y) &= \sum_{n=1}^N E_{D_n} [\delta(y - D_n)] \\ &= \sum_{n=1}^N P_{D_n}(y) \\ &= N \frac{\alpha(y)}{\alpha_0} . \end{aligned} \quad (\text{A.19})$$

The correlation function is

$$\begin{aligned}
E_{\bar{n}} [\bar{n}(y)\bar{n}(y')] &= \sum_{n=1}^N E_{D_n} [\delta(y - D_n)] \\
&= \sum_{n=1}^N \sum_{n'=1}^N E_{D_n, D_{n'}} [\delta(y - D_n) \delta(y - D_{n'})] \\
&= \sum_n E_{D_n} [\delta(y - D_n) \delta(y' - D_n)] + \\
&\quad \sum_{n \neq n'} E_{D_n, D_{n'}} [\delta(y - D_n) \delta(y' - D_{n'})] \\
&= \sum_n \int_{\mathcal{Y}} \frac{\alpha(\tilde{y})}{\alpha_0} \delta(y - \tilde{y}) \delta(y' - \tilde{y}) + \\
&\quad \sum_{n \neq n'} \int_{\mathcal{Y}} \int_{\mathcal{Y}} \frac{\alpha(\tilde{y})\alpha(\tilde{y}') + \alpha(\tilde{y})\delta(\tilde{y} - \tilde{y}')}{\alpha_0(\alpha_0 + 1)} \delta(y - \tilde{y}) \delta(y' - \tilde{y}') \\
&= N \frac{\alpha(y)}{\alpha_0} \delta(y - y') + N(N-1) \frac{\alpha(y)\alpha(y') + \alpha(y)\delta(y - y')}{\alpha_0(\alpha_0 + 1)} \\
&= \frac{N}{\alpha_0(\alpha_0 + 1)} [(N-1)\alpha(y)\alpha(y') + (\alpha_0 + N)\alpha(y)\delta(y - y')] .
\end{aligned} \tag{A.20}$$

#### A.6.4 Continuous aggregation

If  $\bar{n}$  is a Dirichlet-Multinomial process over a Euclidean space  $\mathcal{Y}$ , then conditioning on its discrete aggregation  $n'$  produces independent Dirichlet-Multinomial processes  $\{\bar{n}(y) : y \in \mathcal{S}(z)\} \sim \text{DMP}\left(n'(z), \{\alpha(y) : y \in \mathcal{S}(z)\}\right)$  over the partition spaces  $\mathcal{S}(z)$ .

The previous result can be extended to conditioning on a continuous aggregation. Define  $\bar{n} \sim \text{DMP}(N, \alpha)$  over the set  $\mathcal{Y} \times \mathcal{X}$  and the aggregation  $\text{DMP } n' = \int_{\mathcal{Y}} \bar{n}(y, \cdot) dy$  over set  $\mathcal{X}$  with parameterizing function  $\alpha' = \int_{\mathcal{Y}} \alpha(y, \cdot) dy$ .

Use the aggregation property to introduce a Dirichlet-Multinomial process  $\tilde{n}(y; k) = \int_{\Delta_k}^{\Delta(k+1)} \bar{n}(y, x) dx$  with parameter  $\tilde{\alpha}(y; k) = \int_{\Delta_k}^{\Delta(k+1)} \alpha(y, x) dx$ . Additionally, introduce its own aggregation, a Dirichlet-Multinomial random process  $\dot{n}(k) = \int_{\mathcal{Y}} \tilde{n}(y, k) dy$  with parameter  $\dot{\alpha}(k) = \int_{\mathcal{Y}} \tilde{\alpha}(y, k) dy$ . By the conditioning property for discrete aggregations demonstrated previously,  $\tilde{n}(\cdot, k) | \dot{n}(k) \sim \text{DMP}(\dot{n}(k), \tilde{\alpha}(\cdot, k))$  are independent DMP's.

Note that as  $\Delta \rightarrow 0$ ,  $\tilde{n}(y, k) \approx \Delta \bar{n}(y, \Delta k)$ ,  $\tilde{\alpha}(y, k) \approx \Delta \alpha(y, \Delta k)$ , and  $\dot{n}(k) \approx \Delta n'(\Delta k)$ . Letting  $x \equiv \Delta k$ , the statistics of the DMP conditioned on its continuous aggregation can be

represented as

$$\Delta \bar{n}(\cdot, x) | \Delta n'(k) \sim \text{DMP} (\Delta n'(k), \Delta \alpha(\cdot, x)) . \quad (\text{A.21})$$

# Appendix B

PGR: Many sections redundant given Dirichlet perspective...

PGR: needs notation/formatting scrub

## B.1 Maximum *a Posteriori* estimate of $\theta$ given $D$

PGR: delete? been done...

PGR: REDO USING EASY LOG-LIKELIHOOD and LAGRANGE!!!!!! Generalize for Dirichlet!!!

To determine the MAP estimate of the model PMF  $\theta$  given the training data  $D$ ,

$$\hat{\theta}_{MAP}(D) = \arg \max_{\theta \in \Theta} P_{\theta|D}(\theta|D), \quad (\text{B.1})$$

we perform constrained optimization. Note that the set  $\Theta = \left\{ \theta \in \mathbb{R}_{\geq 0}^M : \sum_{m=1}^M \theta_m = 1 \right\}$  implies both equality and inequality constraints...

# Bibliography

- [1] James O. Berger. *Statistical Decision Theory and Bayesian Analysis*. Second. Springer, 1980.
- [2] José M. Bernardo and Adrian F. M. Smith. *Bayesian Theory*. Wiley Series in Probability and Statistics. John Wiley & Sons, 2000.
- [3] Christopher M. Bishop. *Pattern Recognition and Machine Learning*. Information Science and Statistics. Springer, 2006.
- [4] George E.P. Box and George C. Tiao. *Bayesian Inference in Statistical Analysis*. Wiley, 1973.
- [5] Richard A. Brualdi. *Introductory Combinatorics*. Fifth. Pearson, 2010.
- [6] William Feller. *An Introduction to Probability Theory and Its Applications*. Second. Vol. 2. Probability and Mathematical Statistics. New York, New York: John Wiley & Sons, 1971.
- [7] Thomas S. Ferguson. “A Bayesian Analysis of Some Nonparametric Problems”. In: *The Annals of Statistics* 1.2 (1973), pp. 209–230.
- [8] Ronald L. Graham, Donald E. Knuth, and Oren Patashnik. *Concrete Mathematics*. Second. Reading, Massachusetts: Addison-Wesley, 1994.
- [9] Norman L. Johnson, Samuel Kotz, and N. Balakrishnan. *Discrete Multivariate Distributions*. Probability and Statistics. John Wiley & Sons, 1997.
- [10] Steven M. Kay. *Fundamentals of Statistical Signal Processing: Detection Theory*. Vol. 2. Signal Processing Series. Upper Saddle River, New Jersey: Prentice-Hall, 1998.

- [11] Thomas P. Minka. *Bayesian inference, entropy, and the multinomial distribution*. Tech. rep. Microsoft Research, 2003.
- [12] Kevin P. Murphy. *Binomial and multinomial distributions*. Tech. rep. University of British Columbia, 2006.
- [13] Athanasios Papoulis and S. Unnikrishna Pillai. *Probability, Random Variables and Stochastic Processes*. Fourth. McGraw-Hill, 2002.
- [14] C. Radhakrishna Rao. “Maximum Likelihood Estimation for the Multinomial Distribution”. In: *Sankhyā: The Indian Journal of Statistics (1933-1960)* 18.1/2 (1957), pp. 139–148.
- [15] Steven Roman. “The Logarithmic Binomial Formula”. In: *American Mathematical Monthly* 99.7 (1992).
- [16] Walter Rudin. *Real and Complex Analysis, 3rd Ed.* USA: McGraw-Hill, Inc., 1987. ISBN: 0070542341.
- [17] Frederick F. Stephan. “The Expected Value and Variance of the Reciprocal and other Negative Powers of a Positive Bernoullian Variate”. In: *The Annals of Mathematical Statistics* 16.1 (1945), pp. 50–61.
- [18] Sergios Theodoridis. *Machine Learning: A Bayesian and Optimization Perspective*. Elsevier, 2015.
- [19] J. G. Wendel. “Note on the Gamma Function”. In: *The American Mathematical Monthly* 55.9 (1948), pp. 563–564.