
Gradient Estimation For Exactly- k Constraints

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 The exactly- k constraint is ubiquitous in machine learning and scientific applica-
2 tions, such as ensuring that the sum of electric charges in a neutral atom is zero.
3 However, enforcing such constraints in machine learning models while allowing
4 differentiable learning is challenging. In this work, we aim to provide a “cookbook”
5 for seamlessly incorporating exactly- k constraints into machine learning models
6 by extending a recent gradient estimator from Bernoulli variables to Gaussian and
7 Poisson variables, utilizing constraint probabilities. We show the effectiveness of
8 our proposed gradient estimators in synthetic experiments, and further demonstrate
9 the practical utility of our approach by training neural networks to predict partial
10 charges for metal-organic frameworks, aiding virtual screening in chemistry. Our
11 proposed method not only enhances the capability of learning models but also
12 expands their applicability to a wider range of scientific domains where satisfaction
13 of constraints is crucial.

14 1 Introduction

15 The exactly- k constraint, that is, the sum of n variables is equal to k , is not only ubiquitous in
16 machine learning such as learning sparse features [Chen et al., 2018] and discrete variational auto-
17 encoders [Rolfe, 2016], but also critical to scientific applications such as charge-neutral scenarios in
18 computational chemistry [Raza et al., 2020] and count-aware cell type deconvolution [Liu et al., 2023].
19 In the former cases, the variables are binary while in the latter cases, the variables are continuous
20 or integer, depending on the applications. Such tasks can involve optimizing the expectation of an
21 objective function with respect to variables satisfying the exactly- k constraint, whose distributions
22 are parameterized by neural networks. This optimization problem is challenging since the expectation
23 can be intractable and thus gradient estimation is required. Existing estimators include score-function-
24 based ones that suffer from high variance and reparameterization-based ones that require relaxation
25 and can be highly biased Xie and Ermon [2019]. A recently proposed gradient estimator [Ahmed
26 et al., 2023] outperforms the aforementioned estimators by leveraging constraint probability and
27 avoiding relaxations. Still, it is limited to the exactly- k constraint on Bernoulli variables.

28 In this work, we aim to carry out a systematic study of gradient estimation for exactly- k constraints
29 over Bernoulli, Gaussian, and Poisson variables, the three most commonly used distributions in
30 modeling. We show that on the forward pass, the constrained distributions have closed-form rep-
31 resentations, and thus exact sampling from the constrained distribution can be achieved. On the
32 backward pass, we reparameterize the gradient of the loss function with respect to the samples as a
33 function of the expected marginals of the constrained distributions. Further, we find that under certain
34 loss functions, the expected loss under the constrained distribution has a closed-form solution. That
35 is, in such cases, we are able to train models under the exactly- k constraint without any gradient

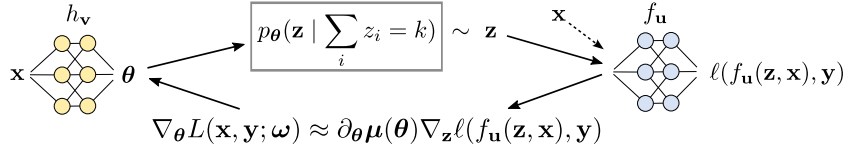


Figure 1: Model formulation under an exactly- k constraint.

36 estimations. We include synthetic experiments to evaluate the bias and variance of our proposed
 37 gradient estimation on Gaussian and Poisson variables. We also include an experiment on predicting
 38 partial charges for metal-organic frameworks, where our gradient estimation, when combined with an
 39 ensemble method, achieves state-of-the-art prediction performance.

40 2 Problem Statement and Motivation

41 We consider models described by the equations

$$\theta = h_v(\mathbf{x}), \quad \mathbf{z} \sim p_\theta(\mathbf{z} \mid \sum_i z_i = k), \quad \hat{\mathbf{y}} = f_u(\mathbf{z}, \mathbf{x}), \quad (1)$$

42 where $\mathbf{x} \in \mathcal{X}$ and $\hat{\mathbf{y}} \in \mathcal{Y}$ denote feature inputs and target outputs, respectively, $h_v : \mathcal{X} \rightarrow \Theta$ and
 43 $f_u : \mathcal{Z} \times \mathcal{X} \rightarrow \mathcal{Y}$ are smooth, parameterized maps. θ are parameters inducing a distribution over
 44 the latent vector \mathbf{z} and the induced distribution $p_\theta(\mathbf{z})$ is defined as $p_\theta(\mathbf{z}) = \prod_{i=1}^n p_{\theta_i}(z_i)$, with
 45 $p_{\theta_i}(z_i)$ as defined in Table 1, where $\mathcal{N}(z; \mu, \sigma^2)$ denotes the density of a Gaussian distribution with
 46 mean μ and variance σ^2 at z . An exactly- k constraint is enforced over the distribution $p_\theta(\mathbf{z})$,
 47 inducing a conditional distribution $p_\theta(\mathbf{z} \mid \sum_i z_i = k) := p_\theta(\mathbf{z}) \cdot \mathbb{I}[\sum_i z_i = k] / p_\theta(\sum_i z_i = k)$
 48 where the denominator denotes the
 49 constraint probability $p_\theta(\sum_i z_i = k)$.
 50 This formulation is general and it can
 51 subsume neural network models that inte-
 52 grate the exactly- k constraint in the
 53 input, output, or latent space, which we
 54 visualize in Figure 1.

Table 1: Parameterization of the three distribution settings.

VARIABLE	PARAMETERIZED DISTRIBUTION
Bernoulli	$p_{\theta_i}(z_i = 1) = \text{sigmoid}(\theta_i)$ $p_{\theta_i}(z_i = 0) = 1 - \text{sigmoid}(\theta_i)$
Gaussian	$p_{\theta_i}(z_i) = \mathcal{N}(z_i; \mu_i, \sigma_i^2)$ with $\theta_i = (\mu_i, \sigma_i)$
Poisson	$p_{\theta_i}(z_i) = \theta_i^{z_i} e^{-\theta_i} / z_i!$

55 The training of such models is per-

56 formed by optimizing an expected loss to learn parameters $\omega = (\mathbf{v}, \mathbf{u})$ in Equation 1 as below,

$$L(\mathbf{x}, \mathbf{y}; \omega) = \mathbb{E}_{\mathbf{z} \sim p_\theta(\mathbf{z} \mid \sum_i z_i = k)} [\ell(f_u(\mathbf{z}, \mathbf{x}), \mathbf{y})] \quad \text{with } \theta = h_v(\mathbf{x}), \quad (2)$$

57 where $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}^+$ is a point-wise loss function. However, the standard auto-differentiation can
 58 not be directly applied to the expected loss due to two main obstacles. First, for the gradient of L
 59 w.r.t. parameters \mathbf{u} in the decoder network f_u defined as

$$\nabla_{\mathbf{u}} L(\mathbf{x}, \mathbf{y}; \omega) = \mathbb{E}_{\mathbf{z} \sim p_\theta(\mathbf{z} \mid \sum_i z_i = k)} [\partial_{\mathbf{u}} f_u(\mathbf{z}, \mathbf{x})^\top \nabla_{\hat{\mathbf{y}}} \ell(\hat{\mathbf{y}}, \mathbf{y})] \quad (3)$$

60 with $\hat{\mathbf{y}} = f_u(\mathbf{z}, \mathbf{x})$ being decoding of a latent sample \mathbf{z} , the expectation does not allow closed-
 61 form solution in general and requires Monte-Carlo estimations by sampling \mathbf{z} from the constrained
 62 distribution $p_\theta(\mathbf{z} \mid \sum_i z_i = k)$. The same issue arises in the gradient of L w.r.t. parameters \mathbf{v} in the
 63 encoder network defined as

$$\nabla_{\mathbf{v}} L(\mathbf{x}, \mathbf{y}; \omega) = \partial_{\mathbf{v}} h_v(\mathbf{x})^\top \nabla_{\theta} L(\mathbf{x}, \mathbf{y}; \omega). \quad (4)$$

64 The second obstacle lies in the computation of the gradient of L w.r.t. the encoder as in Equation 4
 65 defined as $\nabla_{\theta} L(\mathbf{x}, \mathbf{y}; \omega) := \nabla_{\theta} \mathbb{E}_{\mathbf{z} \sim p_\theta(\mathbf{z} \mid \sum_i z_i = k)} [\ell(f_u(\mathbf{z}, \mathbf{x}), \hat{\mathbf{y}})]$ that requires to compute $\partial_{\theta} \mathbf{z}$,
 66 a derivative that is not well-defined and requires gradient estimation for updating θ . In a recent
 67 work [Ahmed et al., 2023], a gradient estimator called SIMPLE is proposed to tackle these two issues
 68 by *exactly sampling* from the constrained distribution and using *marginals* as a proxy to samples
 69 respectively, where SIMPLE is able to outperform both score-function-based gradient estimators and
 70 reparameterization-based ones. However, SIMPLE is limited to Bernoulli variables and whether the
 71 same gradient estimation can be extended to a larger distribution family remains underexplored.

72 3 Gradient Estimation for Exactly- k

73 We tackle the gradient estimation for the exactly- k constraints by solving the aforementioned two
 74 subproblems: **(P1)** how to sample exactly from the constrained distribution $p_{\theta}(\mathbf{z} \mid \sum_i z_i = k)$ and
 75 **(P2)** how to estimate $\nabla_{\theta} L(\mathbf{x}, \mathbf{y}; \omega)$. By combining solutions to these two problems, we manage to
 76 train the constrained model in an end-to-end manner. Table 3 in the Appendix presents a summary of
 77 the key components in the proposed gradient estimation.

78 3.1 Exact Sampling

79 For both Gaussian and Poisson variables, we find that their constrained distributions conform to
 80 commonly seen closed-form distributions and thus allow efficient sampling by using built-in sampling
 81 algorithms in deep learning frameworks. We formally state our findings below.

82 **Proposition 1** (Gaussian Constrained Distribution). *Given $\mathbf{z} = (z_1, \dots, z_n)^T$ with $z_i \sim \mathcal{N}(\mu_i, \sigma_i^2)$,
 83 the constrained distribution $p(\mathbf{z} \mid \sum_{j=1}^n z_j = k)$ is equivalent to an $n - 1$ dimensional multivariate
 84 normal distribution with mean $\bar{\boldsymbol{\mu}} \in \mathbb{R}^{n-1}$ and covariance matrix $\bar{\boldsymbol{\Sigma}} \in \mathbb{R}^{(n-1) \times (n-1)}$ with their
 85 entries defined as below,*

$$\bar{\boldsymbol{\mu}}_i = \sum_{j=1}^{n-1} \left(\mathbb{1}[i=j] \sigma_i^2 - \frac{\sigma_i^2 \sigma_j^2}{\sum_{i=1}^n \sigma_i^2} \right) \left(c + \frac{\mu_j}{\sigma_j^2} \right) \text{ and } \bar{\boldsymbol{\Sigma}}_{i,j} = \begin{cases} \sigma_i^2 - \frac{(\sigma_i^2)^2}{\sum_{i=1}^n \sigma_i^2} & i=j \\ -\frac{\sigma_i^2 \sigma_j^2}{\sum_{i=1}^n \sigma_i^2} & i \neq j \end{cases}.$$

86 **Proposition 2** (Poisson Constrained Distribution). *Given $\mathbf{z} = (z_1, \dots, z_n)^T$ with $z_i \sim \text{Poisson}(\theta_i)$,
 87 the constrained distribution $p(\mathbf{z} \mid \sum_{j=1}^n z_j = k)$ is equivalent to a multinomial distribution with
 88 parameter k and probabilities $\frac{\theta_1}{\sum_{j=1}^n \theta_j}, \dots, \frac{\theta_n}{\sum_{j=1}^n \theta_j}$.*

89 3.2 Conditional Marginals as Proxy

90 For estimating gradient $\nabla_{\theta} L(\mathbf{x}, \mathbf{y}; \omega)$, we follow an approximation adopted by Ahmed et al. [2023],
 91 Niepert et al. [2021] where the main intuition is to use the conditional marginals $\boldsymbol{\mu} := \mu(\boldsymbol{\theta}) :=$
 92 $\{p_{\theta}(z_j \mid \sum_i z_i = k)\}_{j=1}^n$ as a proxy for samples \mathbf{z} , that is,

$$\nabla_{\theta} L(\mathbf{x}, \mathbf{y}; \omega) \approx \partial_{\theta} \boldsymbol{\mu}(\boldsymbol{\theta}) \nabla_{\mathbf{z}} \ell(\mathbf{x}, \mathbf{y}; \omega), \quad (5)$$

93 where the sample \mathbf{z} is reparameterized to be a function of the conditional marginals and is assumed
 94 to be $\partial_{\boldsymbol{\mu}} \mathbf{z} \approx \mathbf{I}$. In the case of Gaussian and Poisson variables, the reparameterization is achieved
 95 by using the expected marginals conditioning on the exactly- k constraint, that is, $\boldsymbol{\mu} := \mu(\boldsymbol{\theta})$ with
 96 $\boldsymbol{\mu}_j = \mathbb{E}_{p_{\theta}(z_j \mid \sum_i z_i = k)}[z_j]$ as a function of the parameters $\boldsymbol{\theta}$. For succinctness, we refer to $\boldsymbol{\mu}$ as
 97 expected marginals. The remaining question is how to obtain the expected marginals $\boldsymbol{\mu}$. We find that
 98 the expected marginals in both cases have closed-form solutions.

99 **Proposition 3** (Gaussian Conditional Marginal). *Given $\mathbf{z} = (z_1, \dots, z_n)^T$ with $z_i \sim \mathcal{N}(\mu_i, \sigma_i^2)$, the
 100 conditional marginal $p(z_i \mid \sum_{j=1}^n z_j = k)$ follows a univariate Gaussian distribution with mean
 101 $\tilde{\mu}_i = \mu_i + \frac{\sigma_i^2}{\sum_{j=1}^n \sigma_j^2} (k - \sum_{j=1}^n \mu_j)$ and variance $\tilde{\sigma}_i^2 = \sigma_i^2 - \frac{(\sigma_i^2)^2}{\sum_{j=1}^n \sigma_j^2}$, that is, $\boldsymbol{\mu}_i = \tilde{\boldsymbol{\mu}}_i$.*

102 **Proposition 4** (Poisson Conditional Marginal). *Given $\mathbf{z} = (z_1, \dots, z_n)^T$ with $z_i \sim \text{Poisson}(\theta_i)$,
 103 the conditional marginal of $p(z_i \mid \sum_{j=1}^n z_j = k)$ follows a binomial distribution with parameter k
 104 and probability $\frac{\theta_i}{\sum_{j=1}^n \theta_j}$, with $\boldsymbol{\mu}_i = \frac{k \theta_i}{\sum_{j=1}^n \theta_j}$.*

105 3.3 Closed-form Expected Loss

106 This section focuses on some special cases where the expected loss in Equation 2 has a closed-form
 107 solution and thus no gradient estimation is needed. We find that when the decoder $f_{\mathbf{u}}$ is an identity
 108 function, that is, $\mathbf{y} = \mathbf{z}$, the expected loss defined over Gaussian variables has a closed-form solution
 109 when the element-wise loss is L1 loss or L2 loss. The same conclusion holds for Poisson variables
 110 with the element-wise loss being L2 loss. We refer the readers to Proposition 5 and Proposition 6
 111 respectively in Appendix for details.

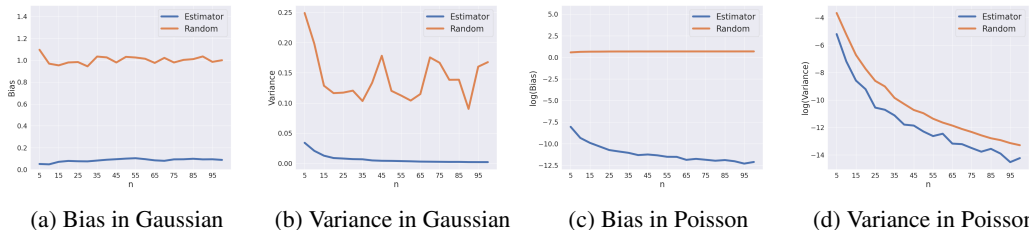


Figure 2: A comparison of our gradient estimation and random estimations on bias and variance.

112 4 Experiments

113 We evaluate our proposed gradient estimation on both synthetic settings and a scientific application.

114 **Synthetic Experiments.** We analyze our proposed gradient estimators for Gaussian and Poisson
 115 variables using three metrics, bias, variance, and averaged error, in synthetic settings where the
 116 ground truth gradients can be obtained by taking derivatives of the closed-form expected loss as
 117 stated in Section 3.3. The distance between the estimated and the ground truth gradient vectors is
 118 measured by the cousin distance defined as $1 - \text{cosine similarity}$. We further compare with a random
 119 estimation as a baseline. Bias and variance results are presented in Figure 2 with additional details
 120 and results presented in Section C in the Appendix, where our proposed gradient estimator is able to
 121 achieve significantly lower bias, variances as well as averaged errors than the baseline, indicating its
 122 effectiveness.

123 **Partial Charge Predictions for Metal-Organic Frameworks.** Metal-organic frameworks (MOFs)
 124 represent a class of materials with a wide range of applications in chemistry and materials science.
 125 Predicting properties of MOFs, such as partial charges on metal ions, is essential for understanding
 126 their reactivity and performance in chemical processes. However, it is challenging due to the complex
 127 interactions between metal ions and ligands and the requirement that the predictions need to satisfy
 128 the charge neutral constraint, that is, an exactly-zero constraint.

129 We adopt the same model as in Raza et al. [2020]
 130 where the charges are assumed to be Gaussian
 131 variables and the element loss is L1 loss, and
 132 address this problem by training the model lever-
 133 aging our observation in Section 3.3 and using
 134 gradients of the expected loss. We further ob-
 135 serve that using an ensemble of such models
 136 gives predictions that also satisfy the charge-
 137 neutral constraint. The prediction performance
 138 of our two proposed approaches is presented in
 139 Table 2, compared with baseline approaches re-
 140 ported by Raza et al. [2020]. Results show that
 141 training using closed-form expected loss achieves
 142 the same performance as MPNN(variance) which
 143 is considered to be the strongest baseline approach, and when further combined with the ensemble
 method, our approach achieves significantly better predictions.

METHOD	MAD
(charge neutrality enforcement)	mean (std)
Constant Prediction	0.324 (7e-3)
Element-mean (uniform)	0.154 (2e-3)
Element-mean (variance)	0.153 (2e-3)
MPNN (uniform)	0.026 (8e-4)
MPNN (variance, reproduced)	0.0251 (8e-4)
Closed-form (ours)	0.0251 (6e-4)
Closed-form + Ensemble (ours)	0.0235 (5e-4)

Table 2: Comparison on prediction performance.

144 5 Conclusion

145 In this work, we provide an extensive study on differentiable learning under exactly- k constraints
 146 given various distribution families. We further provide empirical studies of our proposed gradient
 147 estimation on both synthetic experiments and a scientific application.

148 **References**

- 149 Kareem Ahmed, Zhe Zeng, Mathias Niepert, and Guy Van den Broeck. Simple: A gradient estimator
150 for k-subset sampling. In *Proceedings of the International Conference on Learning Representations*
151 (*ICLR*), may 2023.
- 152 Yoann Altmann, Steve McLaughlin, and Nicolas Dobigeon. Sampling from a multivariate gaussian
153 distribution truncated on a simplex: a review. In *2014 IEEE Workshop on Statistical Signal*
154 *Processing (SSP)*, pages 113–116. IEEE, 2014.
- 155 Jianbo Chen, Le Song, Martin Wainwright, and Michael Jordan. Learning to explain: An information-
156 theoretic perspective on model interpretation. In *International conference on machine learning*,
157 pages 883–892. PMLR, 2018.
- 158 Yulai Cong, Bo Chen, and Mingyuan Zhou. Fast simulation of hyperplane-truncated multivariate
159 normal distributions. 2017.
- 160 Aditya Grover, Eric Wang, Aaron Zweig, and Stefano Ermon. Stochastic optimization of sorting
161 networks via continuous relaxations. *arXiv preprint arXiv:1903.08850*, 2019.
- 162 Allan Gut. *An Intermediate Course in Probability*. Springer, 2009.
- 163 William Holt and Duy Nguyen. *Introduction to Bayesian Data Imputation*. SSRN, 2023.
- 164 Eric Jang, Shixiang Gu, and Ben Poole. Categorical reparameterization with gumbel-softmax. *arXiv*
165 *preprint arXiv:1611.01144*, 2016.
- 166 Carolyn Kim, Ashish Sabharwal, and Stefano Ermon. Exact sampling with integer linear programs
167 and random perturbations. In *Proceedings of the AAAI Conference on Artificial Intelligence*,
168 volume 30, 2016.
- 169 Zhiyuan Liu, Dafei Wu, Weiwei Zhai, and Liang Ma. Sonar enables cell type deconvolution with
170 spatially weighted poisson-gamma model for spatial transcriptomics. *Nature Communications*, 14
171 (1):4727, 2023.
- 172 Hassan Maatouk, Xavier Bay, and Didier Rullière. A note on simulating hyperplane-truncated
173 multivariate normal distributions. *Statistics & Probability Letters*, 191:109650, 2022.
- 174 Chris J Maddison, Andriy Mnih, and Yee Whye Teh. The concrete distribution: A continuous
175 relaxation of discrete random variables. *arXiv preprint arXiv:1611.00712*, 2016.
- 176 Kenneth S Miller. On the inverse of the sum of matrices. *Mathematics magazine*, 54(2):67–72, 1981.
- 177 Mathias Niepert, Pasquale Minervini, and Luca Franceschi. Implicit mle: backpropagating through
178 discrete exponential family distributions. *Advances in Neural Information Processing Systems*, 34:
179 14567–14579, 2021.
- 180 Ali Raza, Arni Sturluson, Cory M Simon, and Xiaoli Fern. Message passing neural networks for
181 partial charge assignment to metal–organic frameworks. *The Journal of Physical Chemistry C*,
182 124(35):19070–19082, 2020.
- 183 Jason Tyler Rolfe. Discrete variational autoencoders. *arXiv preprint arXiv:1609.02200*, 2016.
- 184 Sang Michael Xie and Stefano Ermon. Reparameterizable subset sampling via continuous relaxations.
185 In *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence*,
186 *IJCAI-19*, pages 3919–3925. International Joint Conferences on Artificial Intelligence Organization,
187 7 2019. doi: 10.24963/ijcai.2019/544.

VARIABLE	SAMPLING	EXPECTED MARGINALS	EXPECTED LOSS
Bernoulli	Proposition 2 in Ahmed et al. [2023]	Theorem 1 in Ahmed et al. [2023]	—
Gaussian	Proposition 1	Proposition 3	Proposition 5
Poisson	Proposition 2	Proposition 4	Proposition 6

Table 3: Summary of exact sampling, expected marginals, and closed-form expected loss.

188 A Related Work

189 A substantial amount of research has been devoted to estimating gradients for categorical random
190 variables. Maddison et al. [2016] Jang et al. [2016] proposed to refactor the non-differentiable
191 sample from a categorical distribution with a differentiable sample from a novel Gumbel-Softmax
192 distribution, which enables automatic differentiation. This paper investigates a more complex
193 distribution, k -subset distribution. Gradient estimation under exactly- k constraints has been widely
194 studied. Existing methods either employ the score function and straight-through estimator or suggest
195 custom relaxation [Kim et al., 2016, Chen et al., 2018, Grover et al., 2019, Xie and Ermon, 2019].
196 Xie and Ermon [2019] extends the Gumbel-softmax technique to k -subsets, enabling backpropagation
197 for k -subset sampling. However, this comes at the trade-off of introducing some bias in the learning
198 process due to the use of relaxed samples. While score function estimators offer a seemingly simple
199 solution, it is widely acknowledged that they are prone to exhibiting exceedingly high variance.
200 A recently introduced gradient estimator known as SIMPLE [Ahmed et al., 2023] surpasses its
201 predecessors but is constrained to Bernoulli random variables.

202 Extensive research has been conducted on numerical sampling from multivariate normal distributions
203 while adhering to various constraints. Altmann et al. [2014] reviewed classical Gibbs Sampling on
204 a standard simplex (samples are positive and sum to one) and proposed using Hamiltonian Monte
205 Carlo(HMC) methods. Efficient sampling method for multivariate normal distribution truncated by
206 hyperplanes($\mathbf{Ax} = \mathbf{b}$, where $\dim(\mathbf{x}) = N$ and $\text{rank}(\mathbf{A}) = n < N$) were investigated by Maatouk
207 et al. [2022] and Cong et al. [2017]. These studies focus on numerical simulations, whereas our
208 approach aims to derive a closed-form solution for the multivariate normal distribution subject to the
209 exactly- k constraint.

210 B Theoretical Results

211 **Proposition 5** (Closed-form Expected Loss under Gaussian). *Let $\mathbf{z} = (z_1, \dots, z_n)^T$, where $z_i \sim$*
212 *$\mathcal{N}(\mu_i, \sigma_i^2)$. Let $\mathbf{b} = (b_1, b_2, \dots, b_n)^T$ be the ground truth vector subject to the equality constraint*
213 *$\sum_{j=1}^n b_j = k$. The L1 loss of \mathbf{z} subject to the constraint $\sum_{j=1}^n z_j = k$ is given by*

$$L(\theta) = \sum_{i=1}^n \tilde{\sigma}_i \sqrt{\frac{2}{\pi}} \exp\left(\frac{-(\tilde{\mu}_i - b_i)^2}{2\tilde{\sigma}_i^2}\right) + (\tilde{\mu}_i - b_i) \operatorname{erf}\left(\frac{\tilde{\mu}_i - b_i}{\sqrt{2\tilde{\sigma}_i^2}}\right),$$

214 *where $\tilde{\mu}_i$ and $\tilde{\sigma}_i^2$ are the mean and variance of the conditional marginal of z_i subject to the constraint.*
215 *$\tilde{\mu}_i = \mu_i + \frac{\sigma_i^2}{\sum_{j=1}^n \sigma_j^2} (k - \sum_{j=1}^n \mu_j)$ and $\tilde{\sigma}_i^2 = \sigma_i^2 - \frac{(\sigma_i^2)^2}{\sum_{j=1}^n \sigma_j^2}$. Further, the L2 loss of \mathbf{z} subject to the*
216 *constraint $\sum_{j=1}^n z_j = k$ is given by*

$$L(\theta) = \sum_{i=1}^n \left[\left(\mu_i - \frac{\sigma_i^2 \sum_{j=1}^n \mu_j}{\sum_{j=1}^n \sigma_j^2} \right)^2 + \sigma_i^2 - \frac{(\sigma_i^2)^2}{\sum_{j=1}^n \sigma_j^2} - 2b_i \left(\mu_i - \frac{\sigma_i^2 \sum_{j=1}^n \mu_j}{\sum_{j=1}^n \sigma_j^2} \right) + b_i^2 \right].$$

217 **Proposition 6** (Closed-form Expected Loss under Poisson). *Let $\mathbf{z} = (z_1, \dots, z_n)^T$, where $z_i \sim$*
218 *Poisson(θ_i). Let $\mathbf{b} = (b_1, b_2, \dots, b_n)^T$ be the ground truth vector subject to the equality constraint*

219 $\sum_{j=1}^n b_j = k$. The L2 loss of \mathbf{z} subject to the constraint $\sum_{j=1}^n z_j = k$ is given by

$$L(\theta) = \sum_{i=1}^n \left[k \left(\frac{\theta_i}{\sum_{j=1}^n \theta_j} \right) \left(1 - \frac{\theta_i}{\sum_{j=1}^n \theta_j} \right) + k^2 \left(\frac{\theta_i}{\sum_{j=1}^n \theta_j} \right)^2 - 2k\theta_i \left(\frac{\theta_i}{\sum_{j=1}^n \theta_j} \right) + \theta_i^2 \right].$$

220 C Additional Experiment Results in Synthetic Settings

221 We carried out a series of experiments to analyze the effectiveness of our gradient estimator from
 222 Gaussian and Poisson variables. Our focus lies on three pivotal metrics: bias, variance, and the average
 223 error. Since, we only care about the direction of the gradients, we employed the cosine distance,
 224 namely 1 - cosine similarity, to measure the deviation of our gradient estimators from the ground truth
 225 vector. The ground truth logits, \mathbf{n} , are sampled from $\mathcal{N}(\mathbf{0}, \mathbf{I})$ satisfying the constraint. We plotted the
 226 three metrics against the dimension of \mathbf{z} , namely n , and graphed the standard deviations. For each n ,
 227 we randomly generated 10 sets of parameters and calculated the metrics for each set. Then, we take
 228 average of these 10 repeats and computed their standard deviations. We compare our results with
 229 random guess. The randomly generated gradients are sampled from $\mathcal{N}(\mathbf{0}, \mathbf{I})$.

230 **Gaussian** We use the L1 loss function $L(\theta) = \mathbb{E}_{\mathbf{z} \sim p_{\theta}(\mathbf{z} | \sum_i z_i = 0)} [\|\mathbf{z} - \mathbf{b}\|_1]$. The constraint, k , is
 231 set to 0. We observe that the bias and average error remain relatively stable across various values of
 232 n , with biases hovering around 0.1 and average errors hovering around 0.3. The variance steadily
 233 decreases and converges to a relatively low value. Our estimator outperforms the baseline across all
 dimensions in all three metrics.

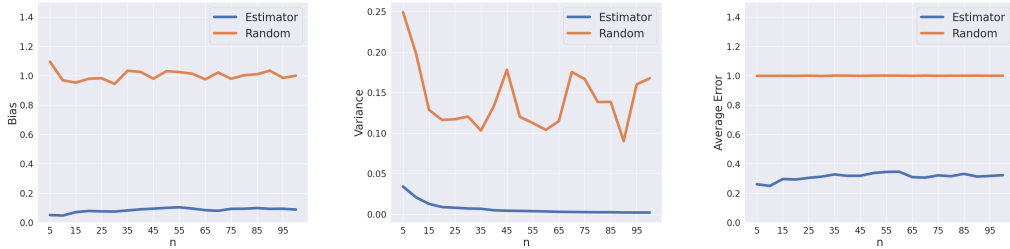


Figure 3: Synthetic Experiment with Gaussian Variables.

234

235 **Poisson** We use the L2 loss function $L(\theta) = \mathbb{E}_{\mathbf{z} \sim p_{\theta}(\mathbf{z} | \sum_i z_i = 0)} [\|\mathbf{z} - \mathbf{b}\|_2^2]$. The constraint is set to
 236 $k = n$. Since, the bias, variance, and average error for our estimators are very small, we opt to take
 their logarithms. In all dimensions and using all three metrics, our estimator surpasses the baseline.

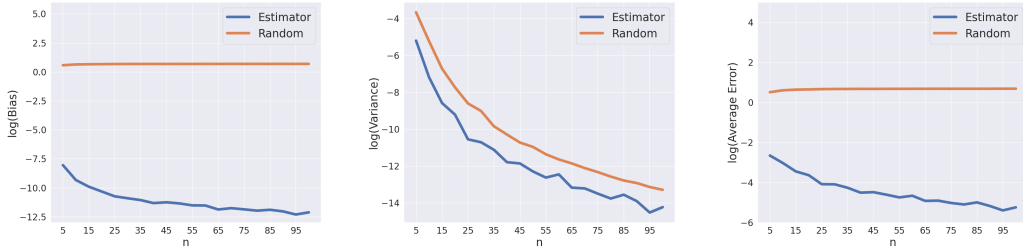


Figure 4: Synthetic Experiment with Poisson Variables.

237

238 **D Additional Experimental Details for Partial Charge Predictions**

239 **Model Architecture** Our model architecture extends the Message Passing Neural Network (MPNN)
 240 Raza et al. [2020] framework and incorporates exact-k constraint for Gaussian variables, ensuring
 241 strict adherence to the critical constraint. The core innovation involves replacing the conventional
 242 L1 loss with the closed-form Gaussian loss function 5. This loss function penalizes deviations
 243 from the exact-k constraint while considering the probabilistic nature of Gaussian variables. This
 244 comprehensive approach not only enables our model to capture complex structural relationships in
 245 MOFs but also ensures accurate predictions of partial charges while respecting the crucial exact-k
 246 constraint, enhancing its applicability in a wide range of graph-based applications, including those
 247 pertaining to metal-organic frameworks.

248 Additionally, we also devise an ensemble methodology to enhance the predictive performance and
 249 robustness of our exact-k constrained MPNN model. To achieve this, we adopt a systematic approach
 250 encompassing model variability, aggregation strategies and cross-validation. Two instances of the
 251 exact-k constrained MPNN model are trained with variations in initialization. We apply the averaging
 252 aggregation technique to combine the predictions from these models. Performance assessment
 253 is conducted through cross-validation techniques. The ensemble’s performance is evaluated on a
 254 separate test dataset to ascertain its generalization ability. This ensemble approach not only elevates
 255 predictive accuracy but also fortifies the model’s resilience, rendering it highly effective for complex
 256 tasks, including those pertaining to metal-organic frameworks.

257 **Training** Here, we describe our training and evaluation process for the exact-k constrained MPNN.
 258 We conducted a random partitioning of the dataset containing 2266 charge-labeled MOFs, creating
 259 distinct training, validation, and test sets (70/10/20%). We use the training set for direct model
 260 parameter tuning, while the validation set aids in hyperparameter selection to prevent overfitting. The
 261 test set plays a crucial role in providing an unbiased assessment of the final model’s performance.

262 **Hyperparameter Tuning** To optimize our model’s performance, we conduct a systematic hyperpa-
 263 rameter tuning process, sequentially optimizing six key hyperparameters: Learning rate, Batch size,
 264 Time steps, Embedding size, Hidden Feature size, and Patience Threshold. After thorough tuning, we
 265 set the hyperparameters to their optimal values: lr = 0.005, batch size = 64, time steps = 4, embedding
 266 size = 20, hidden feature size = 40, and patience threshold = 150, achieving peak model performance.

267 **E Proofs**

268 **E.1 Proposition 1**

269 *Proof.* Let $\mathbf{z} = (z_1, \dots, z_n)^T$, where $z_i \sim \mathcal{N}(\mu_i, \sigma_i^2)$. We attempt to compute a closed-form
 270 solution for the conditional distribution $p\left(\mathbf{z} \mid \sum_{j=1}^n z_j = k\right)$.

$$\begin{aligned} p\left(\mathbf{z} \mid \sum_{j=1}^n z_j = k\right) &= \frac{p\left(\mathbf{z} \cap \sum_{j=1}^n z_j = k\right)}{p\left(\sum_{j=1}^n z_j = k\right)} \\ &= \frac{p(\mathbf{z}) \cdot [\sum_{j=1}^n z_j = k]}{p\left(\sum_{j=1}^n z_j = k\right)} \end{aligned}$$

271 where $[\sum z_i = k]$ is an indicator function. Notice that the denominator $p(\sum_{j=1}^n z_j = k)$ is the
 272 probability distribution function of $Y = \sum_{j=1}^n z_j$ evaluated at k . Since Y is a linear combination of
 273 independent Gaussian random variables, $Y \sim \mathcal{N}(\sum_{j=1}^n \mu_j, \sum_{j=1}^n \sigma_j^2)$. Thus,

$$p\left(\sum_{j=1}^n z_j = k\right) = \frac{1}{\sqrt{2\pi \sum_{j=1}^n \sigma_j^2}} \exp\left[-\frac{1}{2 \sum_{j=1}^n \sigma_j^2} \left(k - \sum_{j=1}^n \mu_j\right)^2\right]$$

274 The joint distribution function $p(\mathbf{z})$, the numerator, follows a multivariate normal distribution with
 275 mean $\mu = (\mu_1, \mu_2, \dots, \mu_n)^T$ and variance $\Sigma = \text{diag}(\sigma_i^2)$. Thus, the conditional distribution can be
 276 rewritten as

$$p\left(\mathbf{z} \mid \sum_{j=1}^n z_j = k\right) = \frac{\prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma_i^2}} \exp\left[-\frac{1}{2\sigma_i^2}(z_i - \mu_i)^2\right]}{\frac{1}{\sqrt{2\pi \sum_{j=1}^n \sigma_j^2}} \exp\left[-\frac{1}{2\sum_{j=1}^n \sigma_j^2}\left(k - \sum_{j=1}^n \mu_j\right)^2\right]} \left[\sum_{j=1}^n z_j = k\right]$$

277 Let $C = \left(\frac{1}{\sqrt{2\pi \sum_{j=1}^n \sigma_j^2}} \exp\left[-\frac{1}{2\sum_{j=1}^n \sigma_j^2}\left(k - \sum_{j=1}^n \mu_j\right)^2\right]\right)^{-1}$. We can express our result as

$$\begin{aligned} p\left(\mathbf{z} \mid \sum_{j=1}^n z_j = k\right) &= C \cdot \left[\sum_{j=1}^n z_j = k\right] \cdot \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma_i^2}} \exp\left[-\frac{1}{2\sigma_i^2}(z_i - \mu_i)^2\right] \\ &= C \cdot f(\mathbf{z}) \end{aligned}$$

278 where $f(\mathbf{z})$ is the joint p.d.f. of the multivariate normal distribution \mathbf{z} . To deal with the indicator
 279 function, let's assume $z_n = k - \sum_{j=1}^{n-1} z_j$. Then, the joint p.d.f. of \mathbf{z} becomes

$$\begin{aligned} f(\mathbf{z}) &= \frac{1}{\sqrt{2\pi\sigma_n^2}} \exp\left[-\frac{1}{2\sigma_n^2}\left(k - \sum_{i=1}^{n-1} z_i - \mu_n\right)^2\right] \cdot \prod_{i=1}^{n-1} \frac{1}{\sqrt{2\pi\sigma_i^2}} \exp\left[-\frac{1}{2\sigma_i^2}(z_i - \mu_i)^2\right] \\ &= (2\pi)^{-\frac{n}{2}} \left(\prod_{i=1}^n \sigma_i\right)^{-1} \\ &\quad \exp\left[-\frac{\left(k^2 - 2k \sum_{i=1}^{n-1} z_i - 2k\mu_n + \left(\sum_{i=1}^{n-1} z_i\right)^2 + 2\mu_n \sum_{i=1}^{n-1} z_i + \mu_n^2 + \sum_{i=1}^{n-1} \frac{z_i^2 - 2z_i\mu_i + \mu_i^2}{\sigma_i^2}\right)}{2}\right] \end{aligned}$$

280 Now, we only consider the terms in the exponential function without $-\frac{1}{2}$.

$$\sum_{i=1}^{n-1} \frac{z_i^2}{\sigma_i^2} + \sum_{i=1}^{n-1} \left(-\frac{2k}{\sigma_n^2} + \frac{2\mu_n}{\sigma_n^2} - \frac{2\mu_i}{\sigma_i^2}\right) z_i + \left(-\frac{2k\mu_n}{\sigma_n^2} + \frac{k^2}{\sigma_n^2} + \frac{\mu_n^2}{\sigma_n^2} + \sum_{i=1}^{n-1} \frac{\mu_i^2}{\sigma_i^2}\right) + \frac{\left(\sum_{i=1}^{n-1} z_i\right)^2}{\sigma_n^2}$$

281 Notice that $\left(\sum_{i=1}^{n-1} z_i\right)^2 = \sum_{i=1}^{n-1} z_i^2 + \sum_{i=1}^{n-1} \sum_{j=1, j \neq i}^{n-1} z_i z_j$. Then, our equation becomes

$$\begin{aligned} &\sum_{i=1}^{n-1} \frac{z_i^2}{\sigma_i^2} + \sum_{i=1}^{n-1} \left(-\frac{2k}{\sigma_n^2} + \frac{2\mu_n}{\sigma_n^2} - \frac{2\mu_i}{\sigma_i^2}\right) z_i + \left(-\frac{2k\mu_n}{\sigma_n^2} + \frac{k^2}{\sigma_n^2} + \frac{\mu_n^2}{\sigma_n^2} + \sum_{i=1}^{n-1} \frac{\mu_i^2}{\sigma_i^2}\right) \\ &\quad + \frac{\sum_{i=1}^{n-1} z_i^2 + \sum_{i=1}^{n-1} \sum_{j=1, j \neq i}^{n-1} z_i z_j}{\sigma_n^2} \\ &= \sum_{i=1}^{n-1} \left(\frac{1}{\sigma_i^2} + \frac{1}{\sigma_n^2}\right) z_i^2 + \sum_{i=1}^{n-1} \left(-\frac{2k}{\sigma_n^2} + \frac{2\mu_n}{\sigma_n^2} - \frac{2\mu_i}{\sigma_i^2}\right) z_i + \left(-\frac{2k\mu_n}{\sigma_n^2} + \frac{k^2}{\sigma_n^2} + \frac{\mu_n^2}{\sigma_n^2} + \sum_{i=1}^{n-1} \frac{\mu_i^2}{\sigma_i^2}\right) \\ &\quad + \frac{\sum_{i=1}^{n-1} \sum_{j=1, j \neq i}^{n-1} z_i z_j}{\sigma_n^2} \\ &= \sum_{i=1}^{n-1} \left[\left(\frac{1}{\sigma_i^2} + \frac{1}{\sigma_n^2}\right) z_i^2 + \frac{\sum_{j=1, j \neq i}^{n-1} z_j}{\sigma_n^2} z_i + \left(-\frac{2k}{\sigma_n^2} + \frac{2\mu_n}{\sigma_n^2} - \frac{2\mu_i}{\sigma_i^2}\right) z_i\right] \\ &\quad + \left(-\frac{2k\mu_n}{\sigma_n^2} + \frac{k^2}{\sigma_n^2} + \frac{\mu_n^2}{\sigma_n^2} + \sum_{i=1}^{n-1} \frac{\mu_i^2}{\sigma_i^2}\right) \end{aligned}$$

282 Then, we consider an arbitrary $n - 1$ dimensional multivariate normal distribution with mean $\bar{\mu}$ and
 283 variance $\bar{\Sigma}$. Its p.d.f. is given by

$$(2\pi)^{-\frac{n-1}{2}} \det \bar{\Sigma}^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\bar{\mathbf{z}} - \bar{\mu})^T \bar{\Sigma}^{-1}(\bar{\mathbf{z}} - \bar{\mu})\right)$$

284 We also only consider the terms in the exponential function without $-\frac{1}{2}$. Let $\bar{\mu}_i$ denotes the i -th
 285 element of the mean $\bar{\mu}$ and $a_{i,j}$ denotes the i,j -th element of the inverse of the variance and covariance
 286 matrix $\bar{\Sigma}^{-1}$.

$$\begin{aligned} &= \bar{\mathbf{z}}^T \bar{\Sigma}^{-1} \bar{\mathbf{z}} - \bar{\mathbf{z}}^T \bar{\Sigma}^{-1} \bar{\boldsymbol{\mu}} - \bar{\boldsymbol{\mu}}^T \bar{\Sigma}^{-1} \bar{\mathbf{z}} + \bar{\boldsymbol{\mu}}^T \bar{\Sigma}^{-1} \bar{\boldsymbol{\mu}} \\ &= \sum_{i=1}^{n-1} \bar{z}_i \left(\sum_{j=1}^{n-1} a_{i,j} \bar{z}_j \right) - \sum_{i=1}^{n-1} \bar{z}_i \left(\sum_{j=1}^{n-1} a_{i,j} \bar{\mu}_j \right) - \sum_{i=1}^{n-1} \bar{\mu}_i \left(\sum_{j=1}^{n-1} a_{i,j} \bar{z}_j \right) + \sum_{i=1}^{n-1} \bar{\mu}_i \left(\sum_{j=1}^{n-1} a_{i,j} \bar{\mu}_j \right) \end{aligned}$$

287 After apply the identity $\sum_{i=1}^{n-1} \bar{z}_i (\sum_{j=1}^{n-1} a_{i,j} \bar{z}_j) = \sum_{i=1}^{n-1} a_{i,i} \bar{z}_i^2 + \sum_{i=1}^{n-1} \sum_{j=1, j \neq i}^{n-1} a_{i,j} \bar{z}_i \bar{z}_j$, the
 288 equation becomes

$$\begin{aligned} &= \sum_{i=1}^{n-1} a_{i,i} \bar{z}_i^2 + \sum_{i=1}^{n-1} \sum_{j=1, j \neq i}^{n-1} a_{i,j} \bar{z}_i \bar{z}_j - \sum_{i=1}^{n-1} \bar{z}_i \left(\sum_{j=1}^{n-1} a_{i,j} \bar{\mu}_j \right) - \sum_{i=1}^{n-1} \bar{\mu}_i \left(\sum_{j=1}^{n-1} a_{i,j} \bar{z}_j \right) \\ &\quad + \sum_{i=1}^{n-1} \bar{\mu}_i \left(\sum_{j=1}^{n-1} a_{i,j} \bar{\mu}_j \right) \\ &= \sum_{i=1}^{n-1} \left[a_{i,i} \bar{z}_i^2 + \bar{z}_i \sum_{j=1, j \neq i}^{n-1} a_{i,j} \bar{z}_j - \left(\sum_{j=1}^{n-1} (a_{i,j} + a_{j,i}) \bar{\mu}_j \right) \bar{z}_i \right] + \sum_{i=1}^{n-1} \bar{\mu}_i \sum_{j=1}^{n-1} a_{i,j} \bar{\mu}_j \end{aligned}$$

289 Now, we consider the terms in the exponent of this arbitrary $n - 1$ dimensional multivariate normal
 290 distribution and the $n - 1$ dimensional multivariate normal distribution we derived previously.

$$\begin{aligned} &\sum_{i=1}^{n-1} \left[\left(\frac{1}{\sigma_i^2} + \frac{1}{\sigma_n^2} \right) z_i^2 + \frac{\sum_{j=1, j \neq i}^{n-1} z_j}{\sigma_n^2} z_i + \left(-\frac{2k}{\sigma_n^2} + \frac{2\mu_n}{\sigma_n^2} - \frac{2\mu_i}{\sigma_i^2} \right) z_i \right] \\ &\quad + \left(-\frac{2k\mu_n}{\sigma_n^2} + \frac{k^2}{\sigma_n^2} + \frac{\mu_n^2}{\sigma_n^2} + \sum_{i=1}^{n-1} \frac{\mu_i^2}{\sigma_i^2} \right) \end{aligned} \quad (6)$$

$$\sum_{i=1}^{n-1} \left[a_{i,i} \bar{z}_i^2 + \bar{z}_i \sum_{j=1, j \neq i}^{n-1} a_{i,j} \bar{z}_j - \left(\sum_{j=1}^{n-1} (a_{i,j} + a_{j,i}) \bar{\mu}_j \right) \bar{z}_i \right] + \sum_{i=1}^{n-1} \bar{\mu}_i \sum_{j=1}^{n-1} a_{i,j} \bar{\mu}_j \quad (7)$$

292 Equation (6) is the term in the exponent of an arbitrary $n - 1$ dimensional multivariate normal
 293 distribution, and Equation (7) is the term in the exponent of previously derived $n - 1$ dimensional
 294 multivariate normal distribution. We get the following three equations by comparing the first few
 295 terms.

$$a_{i,i} = \left(\frac{1}{\sigma_i^2} + \frac{1}{\sigma_n^2} \right) \quad (8)$$

$$a_{i,j} = \frac{1}{\sigma_n^2} \quad (9)$$

$$-\sum_{j=1}^{n-1} (a_{i,j} + a_{j,i}) \bar{\mu}_j = \left(-\frac{2k}{\sigma_n^2} + \frac{2\mu_n}{\sigma_n^2} - \frac{2\mu_i}{\sigma_i^2} \right) \quad (10)$$

296 Equation (8) and (9) define the inverse of the variance and covariance matrix $\bar{\Sigma}^{-1}$. We attempt to
 297 compute $\bar{\Sigma}$. Notice that $\bar{\Sigma}^{-1}$ is equivalent to $\mathbf{A} + \mathbf{B}$, where $\mathbf{A} = \text{diag}(\frac{1}{\sigma_i^2})$ and every element in
 298 matrix \mathbf{B} is $\frac{1}{\sigma_n^2}$.

$$\mathbf{A} = \begin{pmatrix} \frac{1}{\sigma_1^2} & 0 & 0 & \cdots & 0 \\ 0 & \frac{1}{\sigma_2^2} & 0 & \cdots & 0 \\ 0 & 0 & \frac{1}{\sigma_3^2} & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & \frac{1}{\sigma_{n-1}^2} \end{pmatrix} \quad \mathbf{B} = \begin{pmatrix} \frac{1}{\sigma_n^2} & \frac{1}{\sigma_n^2} & \frac{1}{\sigma_n^2} & \cdots & \frac{1}{\sigma_n^2} \\ \frac{1}{\sigma_n^2} & \frac{1}{\sigma_n^2} & \frac{1}{\sigma_n^2} & \cdots & \frac{1}{\sigma_n^2} \\ \frac{1}{\sigma_n^2} & \frac{1}{\sigma_n^2} & \frac{1}{\sigma_n^2} & \cdots & \frac{1}{\sigma_n^2} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \frac{1}{\sigma_n^2} & \frac{1}{\sigma_n^2} & \frac{1}{\sigma_n^2} & \cdots & \frac{1}{\sigma_n^2} \end{pmatrix}$$

299 Consider the following Lemma Miller [1981]

300 **Lemma 1.** Let \mathbf{G} and \mathbf{H} be arbitrary square matrices of the same dimension. If \mathbf{G} and $\mathbf{G} + \mathbf{H}$ are
 301 nonsingular and \mathbf{H} has rank one, then

$$(\mathbf{G} + \mathbf{H})^{-1} = \mathbf{G}^{-1} - \frac{1}{1+g} \mathbf{G}^{-1} \mathbf{H} \mathbf{G}^{-1}$$

302 where $g = \text{tr}(\mathbf{H} \mathbf{G}^{-1})$

303 Since $\det \mathbf{A}$ and $\det(\mathbf{A} + \mathbf{B})$ are nonzero, we know that \mathbf{A} and $\mathbf{A} + \mathbf{B}$ are nonsingular. \mathbf{B} is a rank 1
 304 matrix. By the above lemma, we have

$$(\mathbf{A} + \mathbf{B})^{-1} = \mathbf{A}^{-1} - \frac{1}{1+g} \mathbf{A}^{-1} \mathbf{B} \mathbf{A}^{-1}$$

305 where $g = \text{tr}(\mathbf{B} \mathbf{A}^{-1})$ This is equivalent to

$$\bar{\Sigma} = \mathbf{A}^{-1} - \frac{1}{1 + \text{tr}(\mathbf{B} \mathbf{A}^{-1})} \mathbf{A}^{-1} \mathbf{B} \mathbf{A}^{-1}$$

306 Equation (6) and (7) imply that $\bar{\Sigma}^{-1}$ is a symmetric and positive definite matrix. Its inverse $\bar{\Sigma}$ is also
 307 a symmetric and positive definite matrix. We attempt to find an expression for each element of $\bar{\Sigma}$.
 308 We first consider $\mathbf{B} \mathbf{A}^{-1}$.

$$\begin{aligned} \mathbf{B} \mathbf{A}^{-1} &= \begin{pmatrix} \frac{1}{\sigma_n^2} & \frac{1}{\sigma_n^2} & \frac{1}{\sigma_n^2} & \cdots & \frac{1}{\sigma_n^2} \\ \frac{1}{\sigma_n^2} & \frac{1}{\sigma_n^2} & \frac{1}{\sigma_n^2} & \cdots & \frac{1}{\sigma_n^2} \\ \frac{1}{\sigma_n^2} & \frac{1}{\sigma_n^2} & \frac{1}{\sigma_n^2} & \cdots & \frac{1}{\sigma_n^2} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \frac{1}{\sigma_n^2} & \frac{1}{\sigma_n^2} & \frac{1}{\sigma_n^2} & \cdots & \frac{1}{\sigma_n^2} \end{pmatrix} \begin{pmatrix} \sigma_1^2 & 0 & 0 & \cdots & 0 \\ 0 & \sigma_2^2 & 0 & \cdots & 0 \\ 0 & 0 & \sigma_3^2 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & \sigma_{n-1}^2 \end{pmatrix} \\ &= \begin{pmatrix} \frac{\sigma_1^2}{\sigma_n^2} & \frac{\sigma_2^2}{\sigma_n^2} & \frac{\sigma_3^2}{\sigma_n^2} & \cdots & \frac{\sigma_{n-1}^2}{\sigma_n^2} \\ \frac{\sigma_1^2}{\sigma_n^2} & \frac{\sigma_2^2}{\sigma_n^2} & \frac{\sigma_3^2}{\sigma_n^2} & \cdots & \frac{\sigma_{n-1}^2}{\sigma_n^2} \\ \frac{\sigma_1^2}{\sigma_n^2} & \frac{\sigma_2^2}{\sigma_n^2} & \frac{\sigma_3^2}{\sigma_n^2} & \cdots & \frac{\sigma_{n-1}^2}{\sigma_n^2} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \frac{\sigma_1^2}{\sigma_n^2} & \frac{\sigma_2^2}{\sigma_n^2} & \frac{\sigma_3^2}{\sigma_n^2} & \cdots & \frac{\sigma_{n-1}^2}{\sigma_n^2} \end{pmatrix} \end{aligned}$$

309 Notice that $\text{tr}(\mathbf{B} \mathbf{A}^{-1}) = \sum_{i=1}^{n-1} \frac{\sigma_i^2}{\sigma_n^2}$, so $1 + \text{tr}(\mathbf{B} \mathbf{A}^{-1}) = \sum_{i=1}^n \frac{\sigma_i^2}{\sigma_n^2}$. Then we compute $\mathbf{A}^{-1} \mathbf{B} \mathbf{A}^{-1}$

$$\begin{aligned} \mathbf{A}^{-1} \mathbf{B} \mathbf{A}^{-1} &= \begin{pmatrix} \sigma_1^2 & 0 & 0 & \cdots & 0 \\ 0 & \sigma_2^2 & 0 & \cdots & 0 \\ 0 & 0 & \sigma_3^2 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & \sigma_{n-1}^2 \end{pmatrix} \begin{pmatrix} \frac{\sigma_1^2}{\sigma_n^2} & \frac{\sigma_2^2}{\sigma_n^2} & \frac{\sigma_3^2}{\sigma_n^2} & \cdots & \frac{\sigma_{n-1}^2}{\sigma_n^2} \\ \frac{\sigma_1^2}{\sigma_n^2} & \frac{\sigma_2^2}{\sigma_n^2} & \frac{\sigma_3^2}{\sigma_n^2} & \cdots & \frac{\sigma_{n-1}^2}{\sigma_n^2} \\ \frac{\sigma_1^2}{\sigma_n^2} & \frac{\sigma_2^2}{\sigma_n^2} & \frac{\sigma_3^2}{\sigma_n^2} & \cdots & \frac{\sigma_{n-1}^2}{\sigma_n^2} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \frac{\sigma_1^2}{\sigma_n^2} & \frac{\sigma_2^2}{\sigma_n^2} & \frac{\sigma_3^2}{\sigma_n^2} & \cdots & \frac{\sigma_{n-1}^2}{\sigma_n^2} \end{pmatrix} \\ &= \begin{pmatrix} \frac{(\sigma_1^2)^2}{\sigma_n^2} & \frac{\sigma_2^2 \sigma_1^2}{\sigma_n^2} & \frac{\sigma_3^2 \sigma_1^2}{\sigma_n^2} & \cdots & \frac{\sigma_{n-1}^2 \sigma_1^2}{\sigma_n^2} \\ \frac{\sigma_1^2 \sigma_2^2}{\sigma_n^2} & \frac{(\sigma_2^2)^2}{\sigma_n^2} & \frac{\sigma_3^2 \sigma_2^2}{\sigma_n^2} & \cdots & \frac{\sigma_{n-1}^2 \sigma_2^2}{\sigma_n^2} \\ \frac{\sigma_1^2 \sigma_3^2}{\sigma_n^2} & \frac{\sigma_2^2 \sigma_3^2}{\sigma_n^2} & \frac{(\sigma_3^2)^2}{\sigma_n^2} & \cdots & \frac{\sigma_{n-1}^2 \sigma_3^2}{\sigma_n^2} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \frac{\sigma_1^2 \sigma_{n-1}^2}{\sigma_n^2} & \frac{\sigma_2^2 \sigma_{n-1}^2}{\sigma_n^2} & \frac{\sigma_3^2 \sigma_{n-1}^2}{\sigma_n^2} & \cdots & \frac{(\sigma_{n-1}^2)^2}{\sigma_n^2} \end{pmatrix} \end{aligned}$$

310 The variance and covariance matrix $\bar{\Sigma}$ becomes

$$\begin{pmatrix} \sigma_1^2 & 0 & 0 & \dots & 0 \\ 0 & \sigma_2^2 & 0 & \dots & 0 \\ 0 & 0 & \sigma_3^2 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & \sigma_{n-1}^2 \end{pmatrix} - \frac{1}{\sum_{i=1}^n \sigma_i^2} \begin{pmatrix} \frac{(\sigma_1^2)^2}{\sigma_n^2} & \frac{\sigma_2^2 \sigma_1^2}{\sigma_n^2} & \frac{\sigma_3^2 \sigma_1^2}{\sigma_n^2} & \dots & \frac{\sigma_{n-1}^2 \sigma_1^2}{\sigma_n^2} \\ \frac{\sigma_1^2 \sigma_2^2}{\sigma_n^2} & \frac{(\sigma_2^2)^2}{\sigma_n^2} & \frac{\sigma_3^2 \sigma_2^2}{\sigma_n^2} & \dots & \frac{\sigma_{n-1}^2 \sigma_2^2}{\sigma_n^2} \\ \frac{\sigma_1^2 \sigma_3^2}{\sigma_n^2} & \frac{\sigma_2^2 \sigma_3^2}{\sigma_n^2} & \frac{(\sigma_3^2)^2}{\sigma_n^2} & \dots & \frac{\sigma_{n-1}^2 \sigma_3^2}{\sigma_n^2} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \frac{\sigma_1^2 \sigma_{n-1}^2}{\sigma_n^2} & \frac{\sigma_2^2 \sigma_{n-1}^2}{\sigma_n^2} & \frac{\sigma_3^2 \sigma_{n-1}^2}{\sigma_n^2} & \dots & \frac{(\sigma_{n-1}^2)^2}{\sigma_n^2} \end{pmatrix}$$

311 Thus, we have the following result:

$$\bar{\Sigma}_{i,j} = \begin{cases} \sigma_i^2 - \frac{(\sigma_i^2)^2}{\sum_{i=1}^n \sigma_i^2} & i = j \\ -\frac{\sigma_i^2 \sigma_j^2}{\sum_{i=1}^n \sigma_i^2} & i \neq j \end{cases}$$

312 Next, we derive an expression for $\bar{\mu}$. Since $\bar{\Sigma}^{-1}$ is symmetric, Equation (10) can be transformed into

$$\begin{aligned} -\sum_{j=1}^{n-1} 2a_{i,j}u_j &= \left(-\frac{2k}{\sigma_n^2} + \frac{2\mu_n}{\sigma_n^2} - \frac{2\mu_i}{\sigma_i^2} \right) \\ \sum_{j=1}^{n-1} a_{i,j}u_j &= \left(\frac{k}{\sigma_n^2} + \frac{\mu_i}{\sigma_i^2} - \frac{\mu_n}{\sigma_n^2} \right) \end{aligned}$$

313 This is equivalent to

$$\bar{\Sigma}^{-1}\bar{\mu} = c\mathbf{1} + \mu_{\text{reduced}} \oslash \sigma_{\text{reduced}}$$

314 where $c = \frac{k-\mu_n}{\sigma_n^2}$, $\mu_{\text{reduced}} = (\mu_1, \dots, \mu_{n-1})^T$, $\sigma_{\text{reduced}} = (\sigma_1^2, \dots, \sigma_{n-1}^2)^T$, and \oslash denotes
315 element-wise division of vectors. The mean $\bar{\mu}$ is expressed as

$$\bar{\mu} = \bar{\Sigma}(c\mathbf{1} + \mu_{\text{reduced}} \oslash \sigma_{\text{reduced}}) \quad (11)$$

316 We also attempt to find an element-wise expression for the mean $\bar{\mu}$. Let's define $s_{i,j} = \bar{\Sigma}_{i,j}$. Then we
317 have

$$s_{i,j} = \mathbb{1}[i=j]\sigma_i^2 - \frac{\sigma_i^2 \sigma_j^2}{\sum_{i=1}^n \sigma_i^2}$$

318 From the equation for $\bar{\mu}$, we know that

$$\begin{aligned} \bar{\mu}_i &= \sum_{j=1}^{n-1} s_{i,j} \left(c + \frac{\mu_j}{\sigma_j^2} \right) \\ &= \sum_{j=1}^{n-1} \left(\mathbb{1}[i=j]\sigma_i^2 - \frac{\sigma_i^2 \sigma_j^2}{\sum_{i=1}^n \sigma_i^2} \right) \left(c + \frac{\mu_j}{\sigma_j^2} \right) \end{aligned}$$

319 Finally, we deal with the constant terms in the exponent.

$$-\frac{2k\mu_n}{\sigma_n^2} + \frac{k^2}{\sigma_n^2} + \frac{\mu_n^2}{\sigma_n^2} + \sum_{i=1}^{n-1} \frac{\mu_i^2}{\sigma_i^2} \quad (12)$$

$$\sum_{i=1}^{n-1} \bar{\mu}_i \sum_{j=1}^{n-1} a_{i,j} \bar{\mu}_j \quad (13)$$

320 Equation (12) is the constant term in the exponential function in the probability distribution function
321 derived by taking the cross section of our n dimensional multivariate normal distribution and a hyper-
322 plane. Equation (13) is the constant term in the exponential function in the probability distribution

323 function of an arbitrary $n - 1$ dimensional multivariate normal distribution. The scaling term from
 324 the exponential term is given by

$$\begin{aligned} & -\frac{2k\mu_n}{\sigma_n^2} + \frac{k^2}{\sigma_n^2} + \frac{\mu_n^2}{\sigma_n^2} + \sum_{i=1}^{n-1} \frac{\mu_i^2}{\sigma_i^2} - \sum_{i=1}^{n-1} \bar{\mu}_i \sum_{j=1}^{n-1} a_{i,j} \bar{\mu}_j \\ & = \frac{(\mu_n - k)^2}{\sigma_n^2} + \mathbf{1}^T (\mu_{\text{reduced,squared}} \circ \sigma_{\text{reduced}}) - \bar{\mu}^T \bar{\Sigma}^{-1} \bar{\mu} \end{aligned}$$

325 where $\mu_{\text{reduced,squared}} = (\mu_1^2, \dots, \mu_{n-1}^2)^T$. We define

$$D = \exp \left[-\frac{1}{2} \left(\frac{(\mu_n - k)^2}{\sigma_n^2} + \mathbf{1}^T (\mu_{\text{reduced,squared}} \circ \sigma_{\text{reduced}}) - \bar{\mu}^T \bar{\Sigma}^{-1} \bar{\mu} \right) \right]$$

326 This is our scaling term from the exponent. Finally, we consider the constant term in the front.

$$(2\pi)^{-\frac{n}{2}} \left(\prod_{i=1}^n \sigma_i \right)^{-1} = (2\pi)^{-\frac{1}{2}} \frac{(\prod_{i=1}^n \sigma_i)^{-1}}{\det \bar{\Sigma}^{-\frac{1}{2}}} \cdot (2\pi)^{-\frac{n-1}{2}} \det \bar{\Sigma}^{-\frac{1}{2}}$$

327 $(2\pi)^{-\frac{n}{2}} (\prod_{i=1}^n \sigma_i)^{-1}$ is the constant term of the multivariate normal truncated by the hyperplane, and
 328 $(2\pi)^{-\frac{n-1}{2}} \det \bar{\Sigma}^{-\frac{1}{2}}$ is the constant term of an arbitrary $n - 1$ dimensional multivariate normal. The
 329 scaling term is $E = (2\pi)^{-\frac{1}{2}} \frac{(\prod_{i=1}^n \sigma_i)^{-1}}{\det \bar{\Sigma}^{-\frac{1}{2}}}$. Thus, our conditional distribution is a $n - 1$ dimensional
 330 multivariate normal distribution with p.d.f. given by

$$p \left(\mathbf{z} \mid \sum_{j=1}^n z_j = k \right) = C \cdot D \cdot E \cdot (2\pi)^{-\frac{n-1}{2}} \det \bar{\Sigma}^{-\frac{1}{2}} \exp \left(-\frac{1}{2} (\bar{\mathbf{z}} - \bar{\mu})^T \bar{\Sigma}^{-1} (\bar{\mathbf{z}} - \bar{\mu}) \right)$$

331 where $\bar{\mathbf{z}} = (z_1, \dots, z_{n-1})^T$. □

332 E.2 Proposition 2

333 *Proof.* Let $\mathbf{z} = (z_1, \dots, z_n)^T$, where $z_i \sim \text{Poisson}(\theta_i)$. We attempt to compute a closed-form
 334 solution for the conditional probability $p \left(\mathbf{z} \mid \sum_{j=1}^n z_j = k \right)$.

$$p \left(\mathbf{z} \mid \sum_{j=1}^n z_j = k \right) = \frac{p(\mathbf{z} \cap \sum z_i = k)}{p \left(\sum_{j=1}^n z_j = k \right)}$$

335 Let $Y = \sum_{j=1}^n z_j$. The denominator is the p.d.f. of Y evaluated at k . Since Y is a linear combination
 336 of independent Poisson random variables, we know $Y \sim \text{Poisson}(\sum_{j=1}^n \theta_j)$. Thus,

$$p \left(\sum_{j=1}^n z_j = k \right) = \frac{e^{-\sum_{j=1}^n \theta_j} \left(\sum_{j=1}^n \theta_j \right)^k}{k!}$$

337 Next, let's consider the numerator.

$$p(\mathbf{z} \cap \sum_{j=1}^n z_j = k) = \begin{cases} p(\mathbf{z}) & \sum_{j=1}^n z_j = k \\ 0 & \sum_{j=1}^n z_j \neq k \end{cases}$$

338 where $p(\mathbf{z}) = \prod_{i=1}^n f(z_i) = \prod_{i=1}^n \frac{e^{-\theta_i} \theta_i^{z_i}}{z_i!}$. Thus, our conditional distribution is given by

$$\begin{aligned}
p(\mathbf{z} | \sum_{j=1}^n z_j = k) &= \begin{cases} \frac{e^{-\sum_{i=1}^n \theta_i} \prod_{i=1}^n \theta_i^{z_i}}{\prod_{i=1}^n z_i!} & \sum_{j=1}^n z_j = k \\ 0 & \sum_{j=1}^n z_j \neq k \end{cases} \\
&= \begin{cases} \frac{k! \prod_{i=1}^n \theta_i^{z_i}}{(\sum_{i=1}^n \theta_i)^k \prod_{i=1}^n z_i!} & \sum_{j=1}^n z_j = k \\ 0 & \sum_{j=1}^n z_j \neq k \end{cases} \\
&= \begin{cases} \frac{1}{(\sum_{i=1}^n \theta_i)^k} \cdot \frac{k!}{\prod_{i=1}^n z_i!} \prod_{i=1}^n \theta_i^{z_i} & \sum_{j=1}^n z_j = k \\ 0 & \sum_{j=1}^n z_j \neq k \end{cases} \\
&= \begin{cases} \frac{k!}{\prod_{i=1}^n z_i!} \prod_{i=1}^n \left(\frac{\theta_i}{\sum_{j=1}^n \theta_j} \right)^{z_i} & \sum_{j=1}^n z_j = k \\ 0 & \sum_{j=1}^n z_j \neq k \end{cases} \\
&= f\left(\mathbf{z}; k, \frac{\theta_1}{\sum_{j=1}^n \theta_j}, \dots, \frac{\theta_n}{\sum_{j=1}^n \theta_j}\right)
\end{aligned}$$

339 where $f\left(\mathbf{z}; k, \frac{\theta_1}{\sum_{j=1}^n \theta_j}, \dots, \frac{\theta_n}{\sum_{j=1}^n \theta_j}\right)$ is the probability mass function of a multinomial distribution
340 with parameter k and $\frac{\theta_1}{\sum_{j=1}^n \theta_j}, \dots, \frac{\theta_n}{\sum_{j=1}^n \theta_j}$. \square

341 E.3 Proposition 3

342 *Proof.* Let $\mathbf{z} = (z_1, \dots, z_n)^T$, where $z_i \sim \mathcal{N}(\mu_i, \sigma_i^2)$. We attempt to compute a closed-form solution
343 for the conditional marginal of z_i , $p(z_i | \sum_{j=1}^n z_j = k)$. We first derive the joint distribution of z_i
344 and $\sum_{j=1}^n z_j$. Consider the following affine transformation

$$\mathbf{Az} = \begin{pmatrix} 0 & \dots & 1 & \dots & 0 \\ 1 & \dots & 1 & \dots & 1 \end{pmatrix} \begin{pmatrix} z_1 \\ \vdots \\ z_i \\ \vdots \\ z_n \end{pmatrix} = \begin{pmatrix} z_i \\ \sum_{j=1}^n z_j \end{pmatrix}$$

345 The first row of matrix \mathbf{A} has 1 at i -th column and 0 everywhere, and the last row of matrix \mathbf{A} has 1
346 everywhere.

347 **Theorem 2.** Let $\mathbf{Y} \sim \mathcal{N}_n(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, and let \mathbf{A} be an $m \times n$ matrix of rank m . Then, $\mathbf{AY} \sim$
348 $\mathcal{N}_m(\mathbf{A}\boldsymbol{\mu}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^T)$ Gut [2009]

349 Since matrix \mathbf{A} is full rank, by Theorem 2, $(z_i, \sum_{j=1}^n z_j)^T$ follows a 2 dimensional multivariate
350 normal distribution with mean and variance computed as follows.

$$\begin{aligned}
\mathbf{A}\boldsymbol{\mu} &= \begin{pmatrix} 0 & \dots & 1 & \dots & 0 \\ 1 & \dots & 1 & \dots & 1 \end{pmatrix} \begin{pmatrix} \mu_1 \\ \vdots \\ \mu_i \\ \vdots \\ \mu_n \end{pmatrix} = \begin{pmatrix} \mu_i \\ \sum_{j=1}^n \mu_j \end{pmatrix} \\
\mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^T &= \begin{pmatrix} \sigma_i^2 & & \\ & \sigma_i^2 & \\ & & \sum_{j=1}^n \sigma_j^2 \end{pmatrix}
\end{aligned}$$

351 **Theorem 3.** Suppose that \mathbf{Y} , $\boldsymbol{\mu}$, and $\boldsymbol{\Sigma}$ are partitioned as $\mathbf{Y} = \begin{pmatrix} \mathbf{Y}_1 \\ \mathbf{Y}_2 \end{pmatrix}$, $\boldsymbol{\mu} = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}$, $\boldsymbol{\Sigma} =$
352 $\begin{pmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{pmatrix}$, and $\mathbf{Y} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$. It can be shown that the conditional distribution of \mathbf{Y}_1 given \mathbf{Y}_2

353 is also multivariate normal, $\mathbf{Y}_1 | \mathbf{Y}_2 \sim N(\mu_{1|2}, \Sigma_{1|2})$, where $\mu_{1|2} = \mu_1 + \Sigma_{12}\Sigma_{22}^{-1}(\mathbf{Y}_2 - \mu_2)$,
 354 and $\Sigma_{1|2} = \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21}$ Holt and Nguyen [2023]

355 We apply Theorem 3 to derive the conditional distribution. $z_i | \sum_{j=1}^n z_j \sim \mathcal{N}(\tilde{\mu}_i, \tilde{\sigma}_i^2)$, where the
 356 mean and variance are computed as follows:

$$\tilde{\mu}_i = \mu_i + \frac{\sigma_i^2}{\sum_{j=1}^n \sigma_j^2} (k - \sum_{j=1}^n \mu_j)$$

$$\tilde{\sigma}_i^2 = \sigma_i^2 - \sigma_i^2 \frac{1}{\sum_{j=1}^n \sigma_j^2} \sigma_i^2 = \sigma_i^2 - \frac{(\sigma_i^2)^2}{\sum_{j=1}^n \sigma_j^2}$$

357

□

358 E.4 Proposition 4

359 *Proof.* Let $\mathbf{z} = (z_1, \dots, z_n)^T$, where $z_i \sim \text{Poisson}(\theta_i)$. We attempt to compute a closed-form
 360 solution for the conditional marginal $p(z_i | \sum_{j=1}^n z_j = k)$.

$$p\left(z_i \mid \sum_{j=1}^n z_j = k\right) = \sum \cdots \sum_{(z_1, \dots, z_{i-1}, z_{i+1}, \dots, z_n); \sum_{j=1}^n z_j = k} p(\mathbf{z} \mid \sum z_i = k)$$

$$= \sum \cdots \sum_{(z_1, \dots, z_{i-1}, z_{i+1}, \dots, z_n); \sum_{j=1}^n z_j = k} f\left(\mathbf{z}; k, \frac{\theta_1}{\sum_{j=1}^n \theta_j}, \dots, \frac{\theta_n}{\sum_{j=1}^n \theta_j}\right)$$

361 Since the marginal of each variable of a multinomial distribution is a binomial distribution, then the
 362 conditional marginal is

$$p\left(z_i \mid \sum_{j=1}^n z_j = k\right) = \binom{k}{z_i} \left(\frac{\theta_i}{\sum_{j=1}^n \theta_j}\right)^{z_i} \left(1 - \frac{\theta_i}{\sum_{j=1}^n \theta_j}\right)^{n-z_i}$$

363 This is the probability mass function of a binomial distribution with parameter k and probability
 364 $\frac{\theta_i}{\sum_{j=1}^n \theta_j}$. □

365 E.5 Proposition 5

366 *Proof.* Let $\mathbf{z} = (z_1, \dots, z_n)^T$, where $z_i \sim \mathcal{N}(\mu_i, \sigma_i^2)$. Let $\mathbf{b} = (b_1, b_2, \dots, b_n)^T$ be the ground
 367 truth logits subject to the equality constraint $\sum_{j=1}^n b_j = k$. We attempt to derive a closed-form
 368 solution for the L1 loss of \mathbf{z} subject to the constraint $\sum_{j=1}^n z_j = k$.

$$L(\theta) = \mathbb{E}_{\mathbf{z} \sim p_\theta(\mathbf{z} | \sum_i z_i = k)} [\|\mathbf{z} - \mathbf{b}\|_1]$$

$$= \sum_{i=1}^n \mathbb{E}_{\mathbf{z} \sim p_\theta(\mathbf{z} | \sum_i z_i = k)} [\|z_i - b_i\|_1]$$

369 From previous derivation, we know that the conditional distribution of z_i subject to the equality
 370 constraint is an univariate normal distribution with mean $\tilde{\mu}_i = \mu_i + \frac{\sigma_i^2}{\sum_{j=1}^n \sigma_j^2} (k - \sum_{j=1}^n \mu_j)$ and
 371 variance $\tilde{\sigma}_i^2 = \sigma_i^2 - \frac{(\sigma_i^2)^2}{\sum_{j=1}^n \sigma_j^2}$. Let's define $y_i = z_i - b_i$. Then, $y_i \sim N(\tilde{\mu}_i - b_i, \tilde{\sigma}_i^2)$. Thus,
 372 $\mathbb{E}_{\mathbf{z} \sim p_\theta(\mathbf{z} | \sum_i z_i = k)} [\|y_i\|]$ is the mean of a folded normal distribution.

$$\sum_{i=1}^n \mathbb{E}_{\mathbf{z} \sim p_\theta(\mathbf{z} | \sum_i z_i = k)} [\|y_i\|] = \sum_{i=1}^n \sigma_{y_i} \sqrt{\frac{2}{\pi}} \exp\left(\frac{-\mu_{y_i}^2}{2\sigma_{y_i}^2}\right) + \mu_{y_i} \text{erf}\left(\frac{\mu_{y_i}}{\sqrt{2}\sigma_{y_i}}\right)$$

$$= \sum_{i=1}^n \tilde{\sigma}_i \sqrt{\frac{2}{\pi}} \exp\left(\frac{-(\tilde{\mu}_i - b_i)^2}{2\tilde{\sigma}_i^2}\right) + (\tilde{\mu}_i - b_i) \text{erf}\left(\frac{\tilde{\mu}_i - b_i}{\sqrt{2}\tilde{\sigma}_i}\right)$$

373 We also attempt to derive a closed-form solution for the L2 loss of \mathbf{z} subject to the constraint
 374 $\sum_{j=1}^n z_j = k$.

$$\begin{aligned} L(\theta) &= \mathbb{E}_{\mathbf{z} \sim p_\theta(\mathbf{z} | \sum_i z_i = 0)} [\|\mathbf{z} - \mathbf{b}\|_2^2] \\ &= \sum_{i=1}^n \mathbb{E}_{\mathbf{z} \sim p_\theta(\mathbf{z} | \sum_i z_i = 0)} [z_i^2] - 2 \sum_{i=1}^n \mathbb{E}_{\mathbf{z} \sim p_\theta(\mathbf{z} | \sum_i z_i = 0)} [z_i b_i] + \sum_{i=1}^n \mathbb{E}_{\mathbf{z} \sim p_\theta(\mathbf{z} | \sum_i z_i = 0)} [b_i^2] \end{aligned}$$

375 Since we assume z_i and b_i are independent, and \mathbf{b} is the constant ground truth vector.

$$\begin{aligned} L(\theta) &= \sum_{i=1}^n \mathbb{E}_{\mathbf{z} \sim p_\theta(\mathbf{z} | \sum_i z_i = 0)} [z_i^2] - 2 \sum_{i=1}^n b_i \mathbb{E}_{\mathbf{z} \sim p_\theta(\mathbf{z} | \sum_i z_i = 0)} [z_i] + \sum_{i=1}^n b_i^2 \\ &= \sum_{i=1}^n \mathbb{E}_{z_i \sim p_\theta(z_i | \sum_i z_i = 0)} [z_i^2] - 2 \sum_{i=1}^n b_i \mathbb{E}_{z_i \sim p_\theta(z_i | \sum_i z_i = 0)} [z_i] + \sum_{i=1}^n b_i^2 \end{aligned}$$

376 From previous derivation, we know that the conditional distribution of z_i is $p(z_i | \sum_{j=1}^n z_j = k) =$
 377 $f(z_i; \tilde{\mu}_i = \mu_i + \frac{\sigma_i^2}{\sum_{j=1}^n \sigma_j^2} (k - \sum_{j=1}^n \mu_j), \tilde{\sigma}_i^2 = \sigma_i^2 - \frac{(\sigma_i^2)^2}{\sum_{j=1}^n \sigma_j^2})$. The expectation in the first term is
 378 the second moment of this gaussian distribution.

$$\sum_{i=1}^n \mathbb{E}_{z_i \sim p(z_i | \sum_i z_i = 0)} [z_i^2] = \sum_{i=1}^n \left[\left(\mu_i - \frac{\sigma_i^2 \sum_{j=1}^n \mu_j}{\sum_{j=1}^n \sigma_j^2} \right)^2 + \sigma_i^2 - \frac{(\sigma_i^2)^2}{\sum_{j=1}^n \sigma_j^2} \right]$$

379 Likewise, the expectation in the second term is the mean of this gaussian distribution.

$$\sum_{i=1}^n b_i \mathbb{E}_{z_i \sim p_\theta(z_i | \sum_i z_i = 0)} [z_i] = \sum_{i=1}^n b_i \left(\mu_i - \frac{\sigma_i^2 \sum_{j=1}^n \mu_j}{\sum_{j=1}^n \sigma_j^2} \right)$$

380

□

381 E.6 Proposition 6

382 *Proof.* Let $\mathbf{z} = (z_1, \dots, z_n)^T$, where $z_i \sim \text{Poisson}(\theta_i)$. Let $\mathbf{b} = (b_1, b_2, \dots, b_n)^T$ be the ground
 383 truth vector subject to the equality constraint $\sum_{j=1}^n b_j = k$. We attempt to derive a closed-form
 384 solution for the L2 loss of \mathbf{z} subject to the constraint $\sum_{j=1}^n z_j = k$.

$$\begin{aligned} L(\theta) &= \mathbb{E}_{\mathbf{z} \sim p_\theta(\mathbf{z} | \sum_i z_i = 0)} [\|\mathbf{z} - \mathbf{b}\|_2^2] \\ &= \sum_{i=1}^n \mathbb{E}_{z_i \sim p_\theta(z_i | \sum_j z_j = 0)} [z_i^2] - 2 \sum_{i=1}^n b_i \mathbb{E}_{z_i \sim p_\theta(z_i | \sum_j z_j = 0)} [z_i] + \sum_{i=1}^n b_i^2 \end{aligned}$$

385 Since the conditional marginal distribution is a binomial distribution, it's second moment is given by

$$\sum_{i=1}^n \mathbb{E}_{z_i \sim p_\theta(z_i | \sum_j z_j = 0)} [z_i^2] = \sum_{i=1}^n \left[k \left(\frac{\theta_i}{\sum_{j=1}^n \theta_j} \right) \left(1 - \frac{\theta_i}{\sum_{j=1}^n \theta_j} \right) + k^2 \left(\frac{\theta_i}{\sum_{j=1}^n \theta_j} \right)^2 \right]$$

386 It's first moment(mean) is given by

$$-2 \sum_{i=1}^n b_i \mathbb{E}_{z_i \sim p_\theta(z_i | \sum_j z_j = 0)} [z_i] = -2k \sum_{i=1}^n b_i \left(\frac{\theta_i}{\sum_{j=1}^n \theta_j} \right)$$

387 Thus, we have

$$\begin{aligned} \sum_{i=1}^n \mathbb{E}_{z_i \sim p_\theta(z_i | \sum_j z_j = 0)} [z_i^2] &= \sum_{i=1}^n \left[k \left(\frac{\theta_i}{\sum_{j=1}^n \theta_j} \right) \left(1 - \frac{\theta_i}{\sum_{j=1}^n \theta_j} \right) + k^2 \left(\frac{\theta_i}{\sum_{j=1}^n \theta_j} \right)^2 \right] \\ &\quad - 2k \sum_{i=1}^n b_i \left(\frac{\theta_i}{\sum_{j=1}^n \theta_j} \right) + \sum_{i=1}^n b_i^2 \end{aligned}$$

388

□