

# Master Presentation Marius Hobbhahn

Start: 14:30

## Fast Predictive Uncertainty for Classification with Bayesian Deep Networks

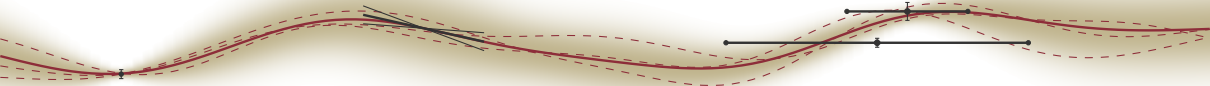
Marius Hobbhahn

30 June 2020

EBERHARD KARLS  
UNIVERSITÄT  
TÜBINGEN



Faculty of Science  
Department of Computer Science  
Chair for the Methods of Machine Learning



# Motivation

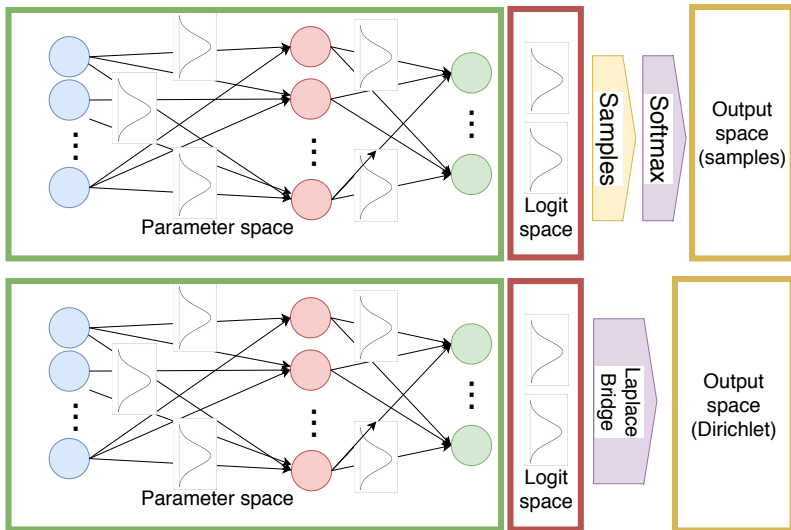
Why do we need fast uncertainty in neural networks?

- ✦ safety-critical applications e.g. self-driving cars
- ✦ trade-off between accuracy and speed
- ✦ out-of-distribution detection



# Context

What's our new contribution?



Theory

# Background

## Change of variable for PDFs

Let  $\mathbf{x}$  be an  $n$ -dimensional continuous random variable with joint density function  $p_{\mathbf{x}}$ . If  $\mathbf{y} = g(\mathbf{x})$ , where  $g$  is a differentiable function, then  $\mathbf{y}$  has density  $p_{\mathbf{y}}$ :

$$p_{\mathbf{y}}(\mathbf{y}) = p_{\mathbf{x}}(g^{-1}(\mathbf{y})) \left| \det \left[ \frac{dg^{-1}(\mathbf{y})}{d\mathbf{y}} \right] \right| \quad (1)$$

where the differential is the Jacobian of the inverse of  $g$  evaluated at  $\mathbf{y}$ .

# A new basis for the Dirichlet

The math

+

$$\text{Dir}(\boldsymbol{\pi}|\boldsymbol{\alpha}) := \frac{1}{B(\boldsymbol{\alpha})} \prod_{k=1}^K \pi_k^{\alpha_k-1} \quad (2)$$

+

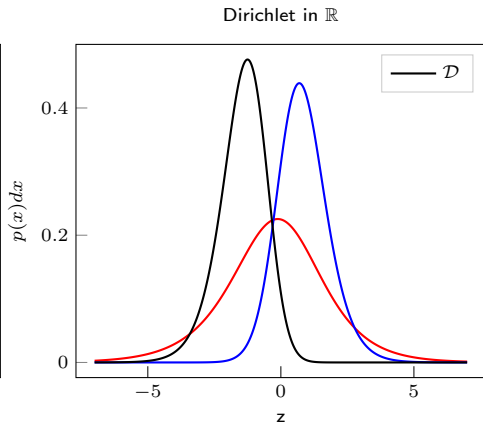
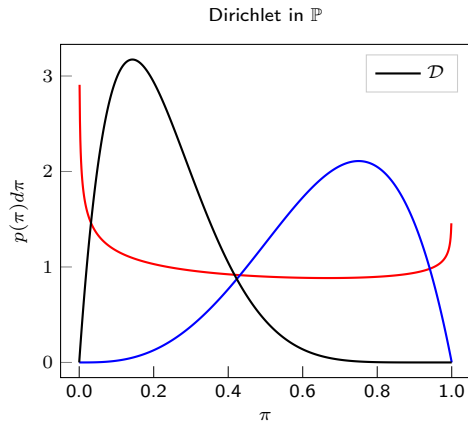
$$\pi_k(\mathbf{z}) := \frac{\exp(z_k)}{\sum_{l=1}^K \exp(z_l)}, \quad (3)$$

+

$$\text{Dir}_{\mathbf{z}}(\boldsymbol{\pi}(\mathbf{z})|\boldsymbol{\alpha}) := \frac{1}{B(\boldsymbol{\alpha})} \prod_{k=1}^K \pi_k(\mathbf{z})^{\alpha_k}, \quad (4)$$

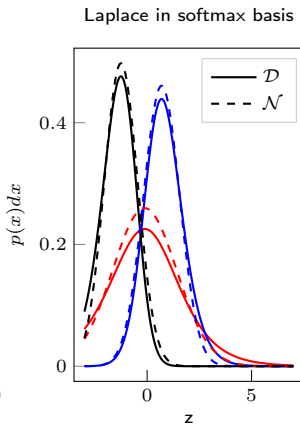
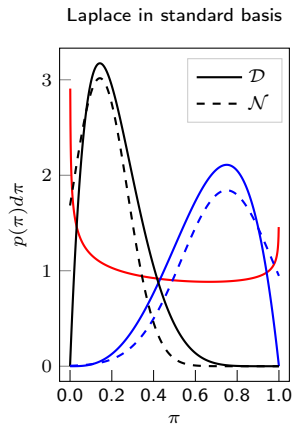
# A new basis for the Dirichlet

In pictures

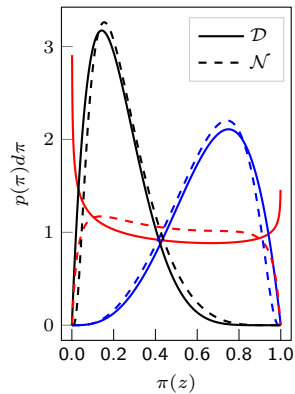


# Building the Bridge

Linking Dirichlet and Gaussian via the Laplace approximation



Transformation back to standard basis





# The Laplace Bridge

A bridge between the parameters of the Dirichlet and Gaussian

$$\alpha_k = \frac{1}{\Sigma_{kk}} \left( 1 - \frac{2}{K} + \frac{e^{\mu_k}}{K^2} \sum_l^K e^{-\mu_l} \right) \quad (5)$$

$$\mu_k = \log \alpha_k - \frac{1}{K} \sum_{l=1}^K \log \alpha_l \quad (6)$$

$$\Sigma_{kl} = \delta_{kl} \frac{1}{\alpha_k} - \frac{1}{K} \left[ \frac{1}{\alpha_k} + \frac{1}{\alpha_l} - \frac{1}{K} \sum_{u=1}^K \frac{1}{\alpha_u} \right] \quad (7)$$

# The Laplace Bridge

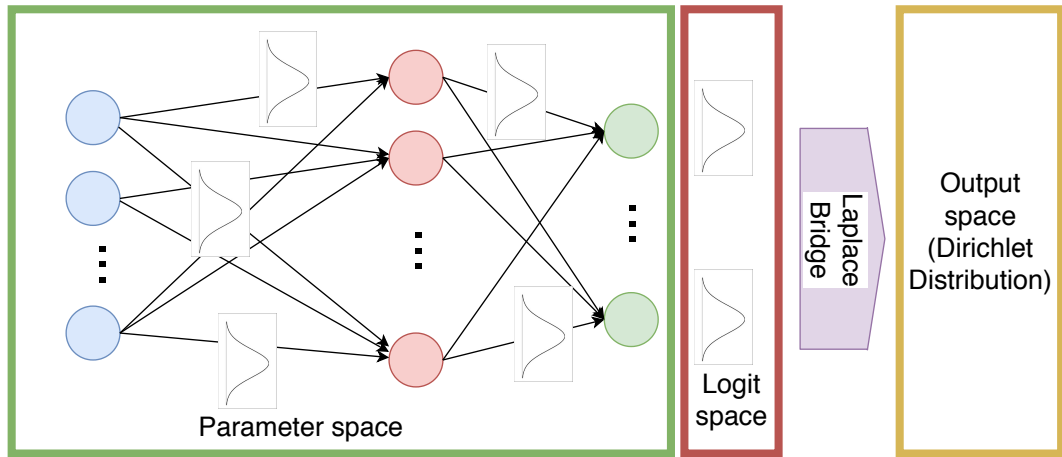
## Summary



- ✦ The Dirichlet in the inverse softmax basis approximates a Gaussian
- ✦ Via the Laplace approximation in the transformed basis we can create a closed-form transformation  $\alpha \rightarrow (\mu, \Sigma)$ .
- ✦ We can also construct an inverse of this transformation  $(\mu, \Sigma) \rightarrow \alpha$
- ✦ In total, we have a **fast** way to transform between the parameters of a Dirichlet and a Gaussian

# The Laplace Bridge

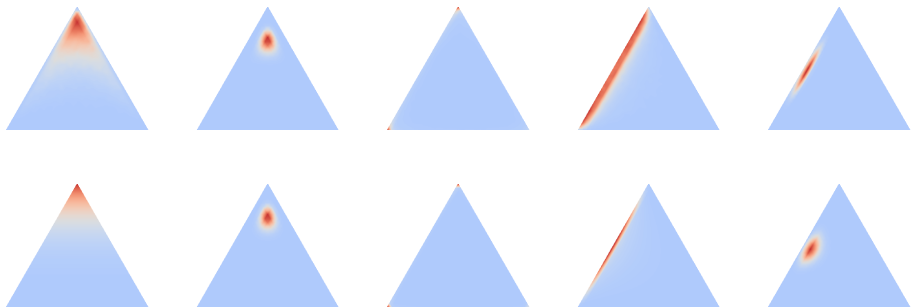
Application to Neural Networks



# Experiments

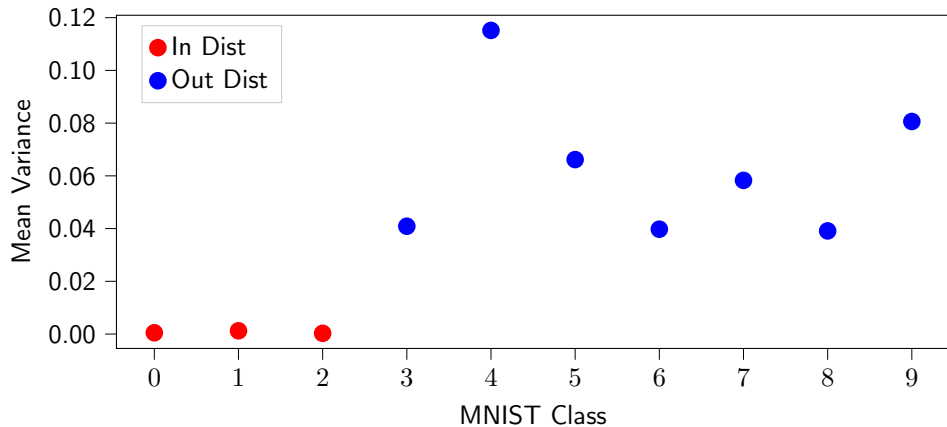
# A sanity check

Samples from a 3D Gaussian + Softmax vs. Dirichlet



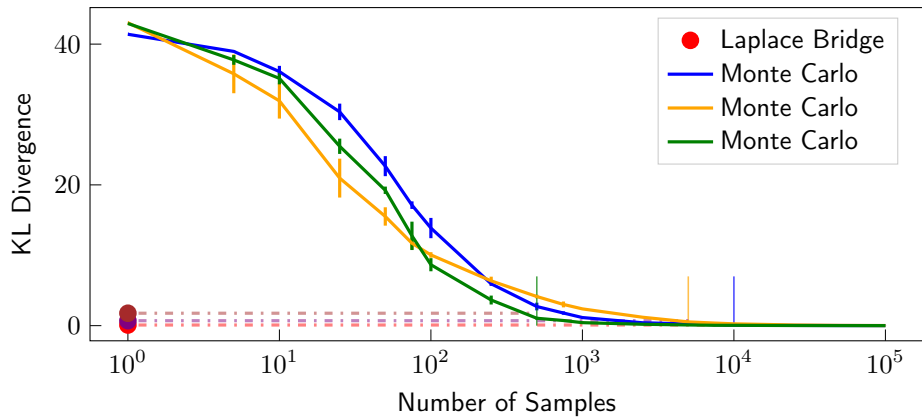
# MNIST

Train on 0,1,2; test on 0-9



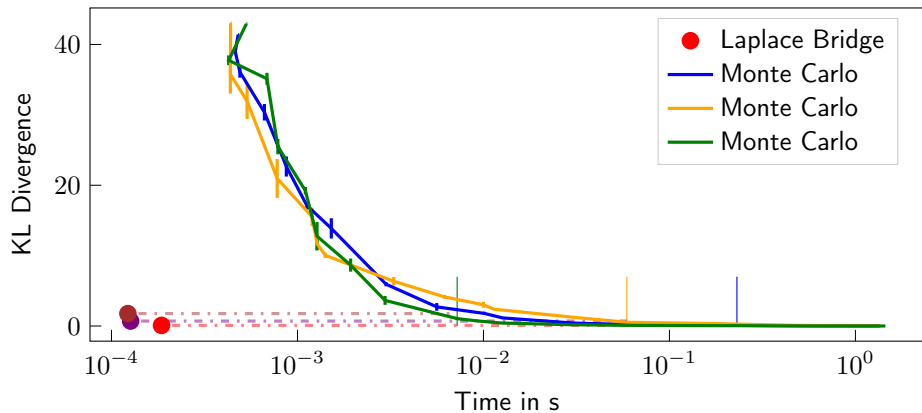
# Speedtest - I

KL divergence vs. number of samples



# Speedtest - II

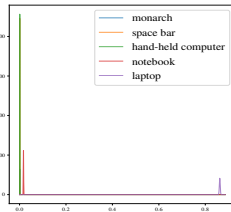
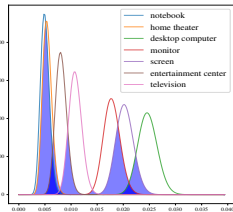
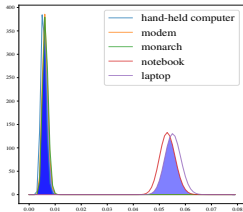
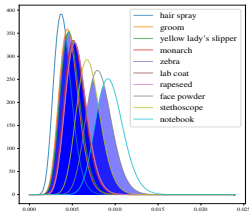
KL divergence vs. wall-clock time





# Imagenet

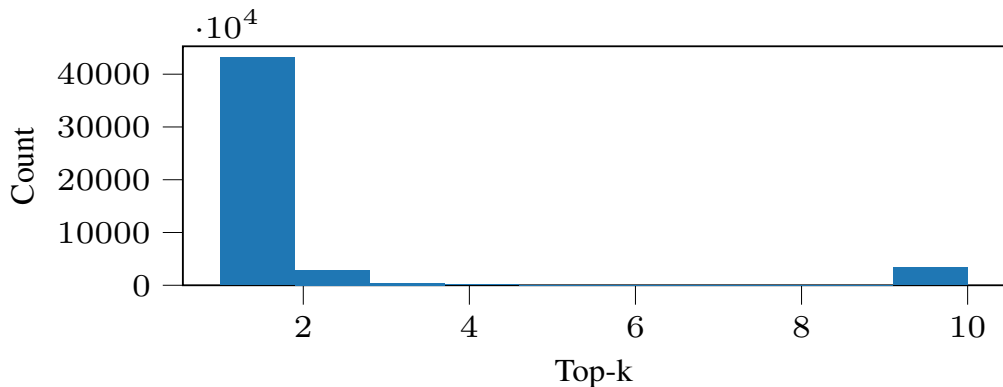
Using the properties of the Dirichlet - The marginal of a Dirichlet is a Dirichlet



We can use the overlap of the distributions to create an uncertainty-aware top-k ranking.

# Imagenet - II

How good is the flexible top-k ranking?



- ✦ The original top-1 accuracy of DenseNet on ImageNet is 0.744 and top-5 accuracy is 0.919
- ✦ The uncertainty-aware top- $k$  accuracy is 0.797, where  $k$  is on average 1.688

# Out-of-distribution Detection

Looking at the numbers

Train	Test	Diag Sampling		Diag LB		KFAC Sampling		KFAC LB		Time in s $\downarrow$	
		MMC $\downarrow$	AUROC $\uparrow$	MMC $\downarrow$	AUROC $\uparrow$	MMC $\downarrow$	AUROC $\uparrow$	MMC $\downarrow$	AUROC $\uparrow$	Sampling	LB
MNIST	MNIST	0.942 $\pm$ 0.007	-	<b>0.987</b> $\pm$ 0.000	-	-	-	-	-	26.8	<b>0.062</b>
MNIST	FMNIST	0.397 $\pm$ 0.001	0.992 $\pm$ 0.000	<b>0.363</b> $\pm$ 0.000	<b>0.996</b> $\pm$ 0.000	-	-	-	-	26.8	<b>0.062</b>
MNIST	notMNIST	<b>0.543</b> $\pm$ 0.000	0.960 $\pm$ 0.000	0.649 $\pm$ 0.000	<b>0.961</b> $\pm$ 0.000	-	-	-	-	50.3	<b>0.117</b>
MNIST	KMNIST	<b>0.513</b> $\pm$ 0.001	<b>0.974</b> $\pm$ 0.000	0.637 $\pm$ 0.000	0.973 $\pm$ 0.000	-	-	-	-	26.9	<b>0.062</b>
CIFAR-10	CIFAR-10	0.948 $\pm$ 0.000	-	<b>0.966</b> $\pm$ 0.000	-	0.857 $\pm$ 0.003	-	<b>0.966</b> $\pm$ 0.000	-	6.58	<b>0.017</b>
CIFAR-10	CIFAR-100	<b>0.708</b> $\pm$ 0.000	<b>0.889</b> $\pm$ 0.000	0.742 $\pm$ 0.000	0.866 $\pm$ 0.000	<b>0.562</b> $\pm$ 0.003	<b>0.880</b> $\pm$ 0.012	0.741 $\pm$ 0.000	0.866 $\pm$ 0.000	6.59	<b>0.016</b>
CIFAR-10	SVHN	<b>0.643</b> $\pm$ 0.000	0.933 $\pm$ 0.000	0.647 $\pm$ 0.000	<b>0.934</b> $\pm$ 0.000	<b>0.484</b> $\pm$ 0.004	<b>0.939</b> $\pm$ 0.001	0.648 $\pm$ 0.003	0.934 $\pm$ 0.001	17.0	<b>0.040</b>
SVHN	SVHN	0.986 $\pm$ 0.000	-	<b>0.993</b> $\pm$ 0.000	-	0.947 $\pm$ 0.002	-	<b>0.993</b> $\pm$ 0.000	-	17.1	<b>0.042</b>
SVHN	CIFAR-100	0.595 $\pm$ 0.000	0.984 $\pm$ 0.000	<b>0.526</b> $\pm$ 0.000	<b>0.985</b> $\pm$ 0.000	<b>0.460</b> $\pm$ 0.004	<b>0.986</b> $\pm$ 0.001	0.527 $\pm$ 0.002	0.985 $\pm$ 0.000	6.62	<b>0.017</b>
SVHN	CIFAR-10	0.593 $\pm$ 0.000	0.984 $\pm$ 0.000	<b>0.520</b> $\pm$ 0.000	<b>0.987</b> $\pm$ 0.000	<b>0.458</b> $\pm$ 0.004	0.986 $\pm$ 0.001	0.520 $\pm$ 0.002	<b>0.987</b> $\pm$ 0.000	6.62	<b>0.017</b>
CIFAR-100	CIFAR-100	<b>0.762</b> $\pm$ 0.000	-	0.590 $\pm$ 0.000	-	0.404 $\pm$ 0.000	-	<b>0.593</b> $\pm$ 0.000	-	6.76	<b>0.016</b>
CIFAR-100	CIFAR-10	0.467 $\pm$ 0.000	0.788 $\pm$ 0.000	<b>0.206</b> $\pm$ 0.000	<b>0.791</b> $\pm$ 0.000	0.213 $\pm$ 0.000	0.788 $\pm$ 0.000	<b>0.209</b> $\pm$ 0.000	<b>0.791</b> $\pm$ 0.000	6.71	<b>0.017</b>
CIFAR-100	SVHN	0.461 $\pm$ 0.000	0.795 $\pm$ 0.000	<b>0.170</b> $\pm$ 0.000	<b>0.815</b> $\pm$ 0.000	0.180 $\pm$ 0.001	<b>0.838</b> $\pm$ 0.001	<b>0.173</b> $\pm$ 0.000	0.815 $\pm$ 0.000	17.3	<b>0.041</b>

- ✦ The Laplace Bridge seems to have better MMC and AUROC compared to sampling from a diagonal Gaussian approximation
- ✦ The Laplace Bridge is as good as a KFAC approximation
- ✦ The Laplace Bridge is around 400 times faster on average

# Conclusions

What can or can't the Laplace Bridge achieve in the context of BNNs?

- ✦ The Laplace Bridge improves an important part of Bayesian Neural Network inference for classification (fast & non-invasive)
- ✦ The Dirichlet distribution has some additional interesting use cases (e.g. the top-k ranking)
- ✦ It will not revolutionize BNNs; it is just one piece in the larger puzzle

Questions?

Future

# The generalized Laplace Bridge

Looking at the larger pattern

- ✦ Similar “Bridges” can be found for all exponential families.
- ✦ Develop a general theoretically grounded framework for the general Laplace Bridge
- ✦ Compute KL-divergences in the different basis

# The generalized Laplace Bridge

So what?

**Implications:** (with a small error)

- ✦ All exponential families can be transformed to Gaussians
- ✦ All exponential families can be transformed to each other
- ✦ All exponential families are conjugate priors for each other



**Backup**

# Backup

## Laplace approximations of a neural network

$$p(c|x) = \mathcal{N}(x; f(x, w_{\text{MAP}}), J(x)^T H^{-1} J(x)) \quad (8)$$

- ✦  $f(x; w_{\text{MAP}})$  is the network output induced by the MAP estimate  $w_{\text{MAP}}$ .
- ✦  $J(x) = \frac{\partial f(x, w_{\text{MAP}})}{\partial w} \in \mathbb{R}^{K \times P}$  is the Jacobian of the network
- ✦  $H_{ij} = \frac{\partial^2 \mathcal{L}(f(x), y)}{\partial w_i \partial w_j} \in \mathbb{R}^{P \times P}$  its Hessian.
- ✦  $K, P$  are the number of classes and parameters of the network respectively.

# Backup

A theoretical bound for the transformation

## Proposition

Let  $\text{Dir}(\boldsymbol{\pi}|\boldsymbol{\alpha})$  be obtained via the Laplace Bridge from a Gaussian distribution  $\mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$  over  $\mathbb{R}^K$ . Then, for each  $k = 1, \dots, K$ , letting  $\alpha_{\neq k} := \sum_{l \neq k} \alpha_l$ , if

$$\alpha_k > \frac{1}{4} \left( \sqrt{9\alpha_{\neq k}^2 + 10\alpha_{\neq k} + 1} - \alpha_{\neq k} - 1 \right),$$

then the variance  $\text{Var}(\pi_k|\boldsymbol{\alpha})$  of the  $k$ -th component of  $\boldsymbol{\pi}$  is increasing in  $\Sigma_{kk}$ .

# Backup

## Computing the Hessian

First, we consider the special case where  $\pi$  is confined to a  $I - 1$  dimensional subspace satisfying  $\sum_i \pi_i = c$ . In this subspace we can represent  $\pi$  by an  $I - 1$  dimensional vector  $\mathbf{a}$  such that

$$\pi_i = a_i \quad i, \dots, I - 1 \quad (9)$$

$$\pi_I = c - \sum_i^{I-1} a_i \quad (10)$$

and similarly we can represent  $\mathbf{z}$  by an  $I - 1$  dimensional vector  $\varrho$ :

$$z_i = \varrho_i \quad i, \dots, I - 1 \quad (11)$$

$$z_I = 1 - \sum_i^{I-1} \varrho_i \quad (12)$$

# Backup

## Computing the Hessian - II

then we can find the density over  $\varrho$  (which is proportional to the required density over  $\mathbf{z}$ ) from the density over  $\boldsymbol{\pi}$  (which is proportional to the given density over  $\boldsymbol{\pi}$ ) by finding the determinant of the  $(I - 1) \times (I - 1)$  Jacobian  $\mathbf{J}$  given by

$$J_{ik} = \frac{\partial \varrho_i}{\partial a_i} = \sum_j^I \frac{\partial z_i}{\partial \pi_j} \frac{\partial \pi_j}{\partial a_k} \quad (13)$$

$$= \delta_{ik} \mathbf{z}_i - \mathbf{z}_i \mathbf{z}_k + \mathbf{z}_i \mathbf{z}_I = \mathbf{z}_i (\delta_{ik} - (\mathbf{z}_k - \mathbf{z}_I)) \quad (14)$$

# Backup

## Computing the Hessian - III

We define two additional  $I - 1$  dimensional helper vectors  $\mathbf{z}_k^+ := \mathbf{z}_k - \mathbf{z}_I$  and  $n_k := 1$ , and use  $\det(I - xy^T) = 1 - x \cdot y$  from linear algebra. It follows that

$$\det J = \prod_{i=1}^{I-1} \mathbf{z}_i \times \det[I - n\mathbf{z}^{+T}] \quad (15)$$

$$= \prod_{i=1}^{I-1} \mathbf{z}_i \times (1 - n \cdot \mathbf{z}^+) \quad (16)$$

$$= \prod_{i=1}^{I-1} \mathbf{z}_i \times \left(1 - \sum_k \mathbf{z}_k^+\right) = I \prod_{i=1}^I \mathbf{z}_i \quad (17)$$