

Information maximising non-linear data compression

Tom Charnock

About me

Current work

- ▶ Started at the IAP in May
- ▶ Statistical problems
- ▶ Novel uses of neural networks
- ▶ IAP machine learning journal club

Previous work

- ▶ Theoretical particle physics PhD at Nottingham
- ▶ Statistical properties of the CMB
- ▶ Cosmic strings
- ▶ Classification of supernovae

Information maximising non-linear data compression

Fisher information

Fisher information matrix

$$\mathbf{F}_{\alpha\beta}(\boldsymbol{\theta}) = - \left\langle \frac{\partial^2 \ln \mathcal{L}(\mathbf{d}|\boldsymbol{\theta})}{\partial \theta_\alpha \partial \theta_\beta} \right\rangle \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_{\text{fid}}}$$

Fisher information of Gaussian likelihood

$$-2 \ln \mathcal{L}(\mathbf{d}|\boldsymbol{\theta}) = (\mathbf{d} - \boldsymbol{\mu}(\boldsymbol{\theta}))^T \mathbf{C}^{-1} (\mathbf{d} - \boldsymbol{\mu}(\boldsymbol{\theta})) + n_{\mathbf{d}} \ln |2\pi\mathbf{C}|$$

$$\mathbf{F}_{\alpha\beta}(\boldsymbol{\theta}) = \frac{1}{2} \text{Tr} \left[\frac{\partial \boldsymbol{\mu}(\boldsymbol{\theta})^T}{\partial \theta_\alpha} \mathbf{C}^{-1} \frac{\partial \boldsymbol{\mu}(\boldsymbol{\theta})}{\partial \theta_\beta} + \frac{\partial \boldsymbol{\mu}(\boldsymbol{\theta})^T}{\partial \theta_\beta} \mathbf{C}^{-1} \frac{\partial \boldsymbol{\mu}(\boldsymbol{\theta})}{\partial \theta_\alpha} \right]$$

Assuming the covariance is independent of the parameters

Information conserving linear data compression

$$n_d \rightarrow n_\theta$$

Linear summaries

$$\tilde{d}_\alpha = \mathbf{r}_\alpha^T \mathbf{d}$$

with each \mathbf{r}_α orthogonal such that \tilde{d}_α are uncorrelated.

$$\mathbf{r}_1 = \frac{\mathbf{C}^{-1} \boldsymbol{\mu}_{,1}}{\sqrt{\boldsymbol{\mu}_{,1}^T \mathbf{C}^{-1} \boldsymbol{\mu}_{,1}}}$$
$$\mathbf{r}_\alpha = \frac{\mathbf{C}^{-1} \boldsymbol{\mu}_{,\alpha} - \sum_{\beta=1}^{\alpha-1} (\boldsymbol{\mu}_{,\alpha}^T \mathbf{r}_\beta) \mathbf{r}_\beta}{\sqrt{\boldsymbol{\mu}_{,1}^T \mathbf{C}^{-1} \boldsymbol{\mu}_{,1} - \sum_{\beta=1}^{\alpha-1} (\boldsymbol{\mu}_{,\alpha}^T \mathbf{r}_\beta)^2}}$$

with $_{,\alpha} \equiv \partial/\partial\theta_\alpha$, $\boldsymbol{\mu} = \boldsymbol{\mu}(\boldsymbol{\theta})$ and \mathbf{C} and $\boldsymbol{\mu}$ evaluated at fiducial parameter values.

Less pleasant likelihoods

Unspecified function of the data $f(\mathbf{d})$

$$-2 \ln \mathcal{L}(\mathbf{d}|\boldsymbol{\theta}) = (\mathbf{f}(\mathbf{d}) - \boldsymbol{\mu}_f)^\top \mathbf{C}_f^{-1} (\mathbf{f}(\mathbf{d}) - \boldsymbol{\mu}_f)$$

with mean and covariance over some simulations \mathbf{s}_i

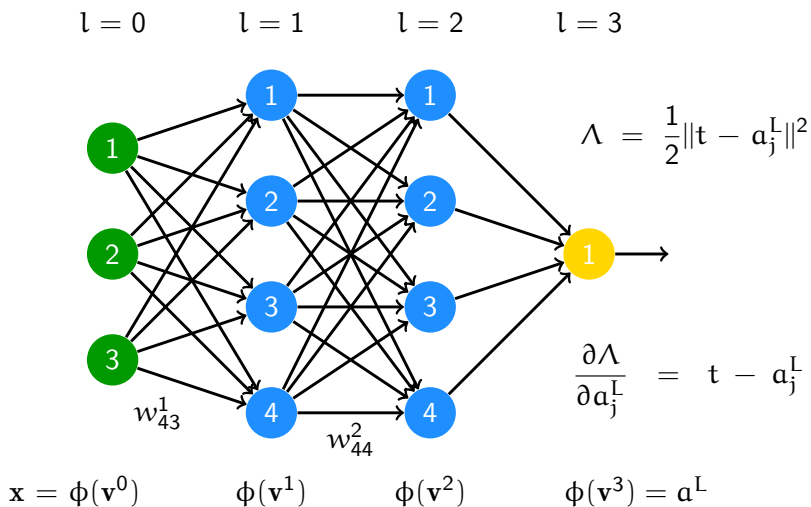
$$\boldsymbol{\mu}_f = \frac{1}{n_s} \sum_{i=1}^{n_s} f(\mathbf{s}_i)$$

$$\mathbf{C}_f = \frac{1}{n_s - 1} \sum_{i=1}^{n_s} (f(\mathbf{s}_i) - \boldsymbol{\mu}_f)(f(\mathbf{s}_i) - \boldsymbol{\mu}_f)^\top$$

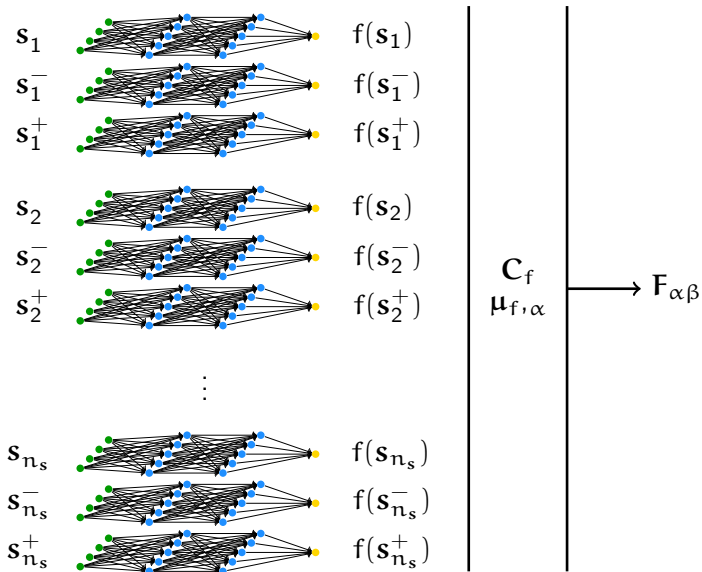
New Fisher information matrix

$$\mathbf{F}_{\alpha\beta} = \frac{1}{2} \text{Tr} \left[\boldsymbol{\mu}_{f,\alpha}^\top \mathbf{C}_f^{-1} \boldsymbol{\mu}_{f,\beta} + \boldsymbol{\mu}_{f,\beta}^\top \mathbf{C}_f^{-1} \boldsymbol{\mu}_{f,\alpha} \right]$$

Neural networks - the unspecified function



Information maximising neural network



Test model - the problem

Unknown variance Gaussian noise

Real data is $\mathbf{d} = \{d_i \sim \mathcal{N}(0, \theta) \mid i = 1 \rightarrow n_d\}$ where θ is unknown and $n_d = 10$.

Analytically the likelihood is

$$-2 \ln \mathcal{L}(\mathbf{d}|\theta) = \frac{1}{\theta} \sum_{i=1}^{n_d} d_i^2 + n_d \ln [2\pi\theta]$$

Single sufficient statistic

$$T = \sum_{i=1}^{n_d} d_i^2$$

The Fisher information is

$$\mathbf{F} = \frac{n_d}{2\theta_{\text{fid}}}$$

Choosing $\theta_{\text{fid}} = 1$ gives $\mathbf{F} = 5$.

Test model - data compression

Linear data compression

The sufficient statistic for this problem is non-linear, i.e.

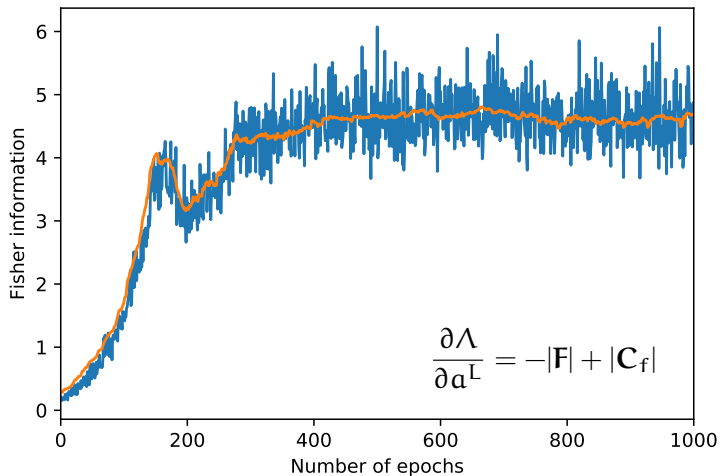
$$T = \sum_{i=1}^{n_d} d_i^2$$

Any linear combination of \mathbf{d} will necessarily give F less than optimal.

Non-linear data compression

Can train the network to become the function f which fully specifies the sufficient statistic?

Test model - training



[96, 96], 0.5 dropout, leaky ReLU with $\alpha = 0.1$, $\eta = 0.01$, $\mathbf{b}_{init}^1 = 0.1$ $\mathbf{w}_{init}^1 = \mathcal{N}(0, \sqrt{2}/\kappa^{1-1})$,
1000 simulations, 200 derivatives, 2 batches, 1000 epochs

Artificial approximate Bayesian computation (AABC)

Get summary of real data $f(\mathbf{d})$

Create simulation $\mathbf{s}_i(\theta)$ within prior range of θ

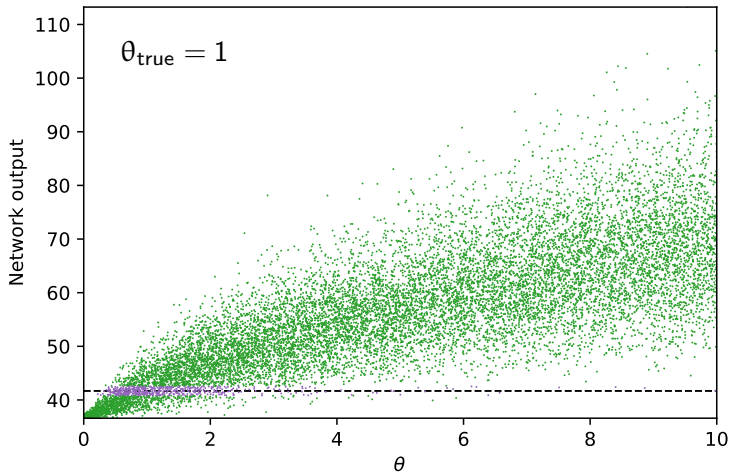
Calculate $f(\mathbf{s}_i(\theta))$

Accept point if

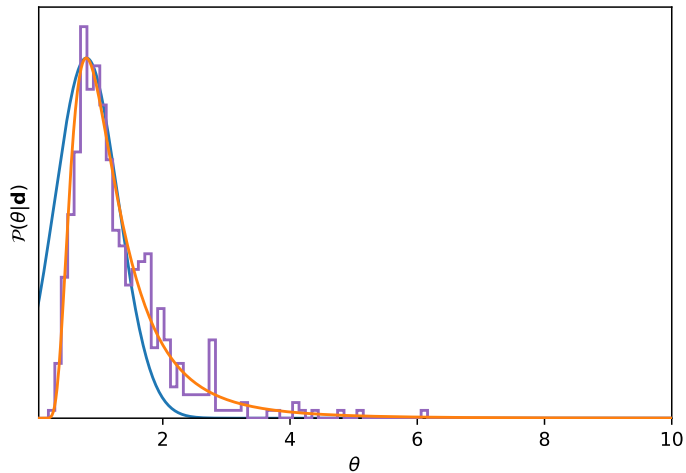
$$\rho = \sqrt{(f(\mathbf{s}_i(\theta)) - f(\mathbf{d}))^2} < \varepsilon$$

Actually use PMC rather than random draws from the prior

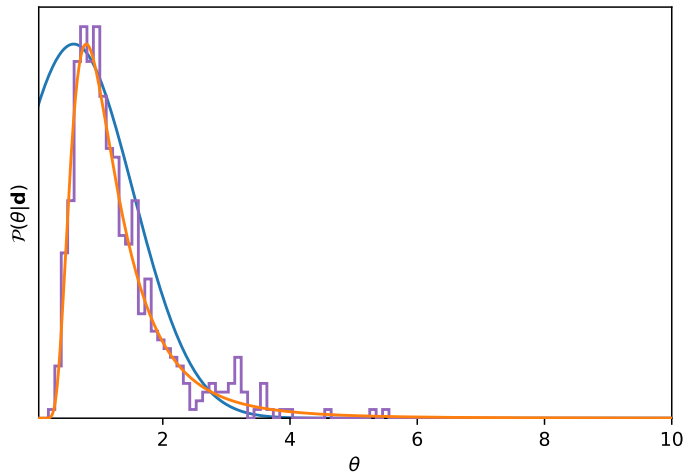
Network output



$$P(\theta|\mathbf{d})$$



$P(\theta|\mathbf{d})$ when $\theta_{\text{fid}} = 2$



Information maximising non-linear data compression

- ▶ Linear data compression can find summaries of Gaussian likelihood (and approximations thereof).
- ▶ Non-linear data compression can find summaries of any data where a function of the data can be written like a Gaussian.
- ▶ The function of the data which maximises the Fisher information is the function which best summarises the data - although this function is unknown.
- ▶ A neural network can be trained (using relatively few training samples) to find the function by maximising the Fisher information as the loss function.
- ▶ This function can now be used for AABC to find $P(\theta|\mathbf{d})$.

About me (non-work-related ice breakers)

Music

- ▶ Avid fan of traditional folk music from all over the world
- ▶ Played in several bands with styles from metal to blues to folk
- ▶ Melodeon
- ▶ Bouzouki
- ▶ Electric guitar

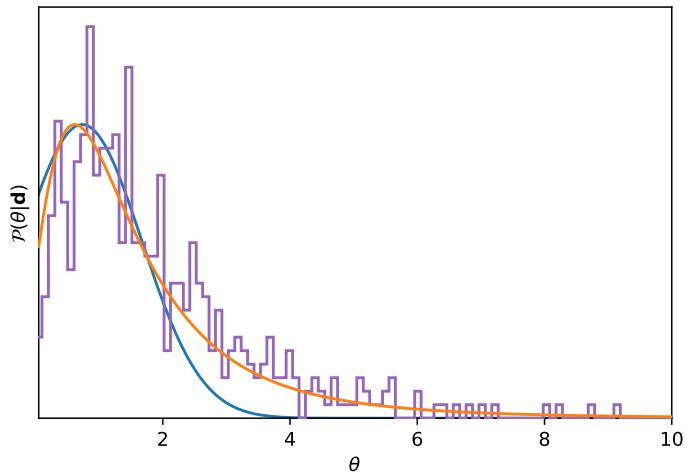


The finer things...

- ▶ A well-made (hipster?) coffee
- ▶ Craft ale
- ▶ Bread and cheese (no wonder I moved to France)
- ▶ Listening to music on vinyl

Etc...

$P(\theta|\mathbf{d})$ when $\theta_{\text{fid}} = 1$ with known noise $\sigma_{\text{noise}}^2 = 1$



$P(\theta|\mathbf{d})$ when $\theta_{\text{fid}} = 1$ with unknown noise $\sigma_{\text{noise}}^2 = [0, 2]$

