

## Aleatoric / Epistemic uncertainty

Aleatoric uncertainty is uncertainty arising from the inherent randomness of an event.  
Unexplainable variance of  $\text{Var}(\theta)$  due to non-deterministic behaviour. Cause for  $\text{Var}(\theta) > 0$ .

Epistemic uncertainty is uncertainty due to lack of knowledge, which in principle can be diminished by increasing the amount of data or with an excessive number of experiments.

## Kullback-Leibler Divergence

$$\begin{aligned} \text{KL}(G, F(\cdot; \theta)) &= \int \log\left(\frac{f(y)}{f_G(y)}\right) dG(y) = \int \log\left(\frac{f(y)}{f(y; \theta)}\right) \cdot g(y) dy \\ &= \underbrace{\int \log\left(\frac{f(y)}{f(y; \theta)}\right) dy}_{\text{empirical KL}} - \int \log(f(y; \theta)) \cdot g(y) dg(y) \\ &= \frac{1}{n} \sum_{i=1}^n \log\left(\frac{g(y_i)}{f(y_i; \theta)}\right) = \frac{1}{n} \sum_{i=1}^n \log(g(y_i)) - \frac{1}{n} \sum_{i=1}^n \log(f(y_i; \theta)) \xrightarrow{n \rightarrow \infty} \text{KL}(G, F(\cdot; \theta)) \end{aligned}$$

To minimize empirical KL we need to maximize  $\sum_i \log(f(y_i; \theta))$

The optimal parameter  $\theta_0$  is defined as  $\theta_0 = \int \frac{\partial \log(f(y; \theta_0))}{\partial \theta} dG(y)$   
and hence estimated through  $\frac{\partial \ell}{\partial \theta}(\hat{\theta}) = 0$

## Properties of Estimates

Bias:  $\text{bias}(\hat{\theta}, \theta_0) = E[\hat{\theta}] - \theta_0$

MSE:  $\text{MSE}(\hat{\theta}, \theta_0) := E[(\hat{\theta} - \theta_0)^2] = \text{Var}(\hat{\theta}) + \text{bias}^2(\hat{\theta}, \theta_0)$

Consistency:  $\text{MSE}(\hat{\theta}, \theta_0) \xrightarrow{n \rightarrow \infty} 0$

Efficiency: We call  $\hat{\theta}$  more efficient than  $\hat{\theta}_0$  if  $\text{MSE}(\hat{\theta}, \theta_0) < \text{MSE}(\hat{\theta}_0, \theta_0)$

Sufficiency:  $t(Y_1, \dots, Y_n)$  is called sufficient for  $\theta$

$P(Y_1 = y_1, \dots, Y_n = y_n | t(Y_1, \dots, Y_n) = t_0; \theta)$   
does not depend on  $\theta$ .

## Neyman-factorisation

$t(Y_1, \dots, Y_n)$  is called sufficient for  $\theta$  iff

$$f(y_1, \dots, y_n; \theta) = h(y_1, \dots, y_n) \cdot g(t(y_1, \dots, y_n), \theta)$$

## Minimal sufficient:

$t(Y_1, \dots, Y_n)$  is minimal sufficient for  $\theta$  if

$t$  is sufficient and for any other sufficient statistic  $t'(Y_1, \dots, Y_n)$

there exists a function  $m$  s.t.  $t(Y_1, \dots, Y_n) = m(t'(Y_1, \dots, Y_n))$

## Example

$Y_i \sim U(0, \pi)$  and  $t(Y_1, \dots, Y_n) = t_0$  and  $n = n \cdot t_0$ .

$$P(Y_1 = y_1, \dots, Y_n = y_n | t(Y_1, \dots, Y_n) = t_0, \pi) = \begin{cases} \frac{1}{(\pi)^n}, & \text{for } \sum_i y_i = n \cdot t_0 \\ 0, & \text{otherwise} \end{cases}$$

Distribution is independent of  $\pi$ .  
 $\Rightarrow t(Y_1, \dots, Y_n) = \bar{y}$  is sufficient.

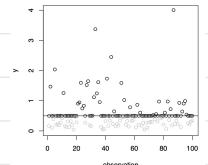
## Likelihood principle

The likelihood principle states that for inference about parameter  $\theta$ , all relevant information is contained in the likelihood function for the data at hand

## Example

Assume  $\log(y) \sim U(\mu, \sigma^2)$ . Instead of  $\min_{\mu} \sum_{i=1}^n (\log(y_i) - \mu)^2$  we solve

$$\max_{\mu} \ell(\mu, \sigma^2) = \max_{\mu} \sum_{i=1}^n I(y_i > c) \cdot \left( -\frac{1}{2} \log(\sigma^2) - \frac{(\log(y_i) - \mu)^2}{2\sigma^2} \right) + \sum_{i=1}^n I(y_i < c) \cdot \Phi\left(\frac{\log(y_i) - \mu}{\sigma}\right)$$



## Estimator Theory

## Maximum Likelihood (MLE)

For a random sample  $y_1, \dots, y_n$  the maximum likelihood estimate is defined as  $\hat{\theta} := \arg \max_{\theta \in \Theta} \ell(\theta; y_1, \dots, y_n)$  which for Fisher-regular distributions occurs when  $s(\theta; y_1, \dots, y_n) = 0$

## Score function

$$s(\theta; y) := \frac{\partial \ell(\theta; y)}{\partial \theta} = \frac{\partial \log f(y; \theta)}{\partial \theta}$$

$$\begin{aligned} \text{Under regularity we have } E_\theta[s(\theta; y)] &= \int \frac{\partial \log f(y; \theta)}{\partial \theta} \cdot f(y; \theta) dy \\ &= \int \frac{\partial f(y; \theta)}{\partial \theta} dy = \frac{\partial}{\partial \theta} \int f(y; \theta) dy = \frac{\partial}{\partial \theta} = 0. \end{aligned}$$

## Fisher Information

$$I(\theta) := \text{Var}_\theta[s(\theta; Y)] = E_\theta[s(\theta; Y)^2] = \int \left( \frac{\partial \log f(y; \theta)}{\partial \theta} \right)^2 f(y; \theta) dy.$$

$$\text{For certain assumptions: } I(\theta) = -E_\theta\left[\frac{\partial^2}{\partial \theta^2} \log f(Y; \theta)\right]$$

$$\text{The observed Fisher information is } J(\theta) = -\sum_{i=1}^n \frac{\partial^2}{\partial \theta^2} \log f(y_i; \theta).$$

? The Fisher information determines how quickly the observed score function converges to the shape of the true score function.

For the case of  $\theta \in \mathbb{R}^n$

$$[I(\theta)]_{ij} := E_\theta\left[\left(\frac{\partial}{\partial \theta_i} \log f(Y; \theta)\right) \left(\frac{\partial}{\partial \theta_j} \log f(Y; \theta)\right)\right]$$

$$\text{For certain assumptions: } [I(\theta)]_{ij} = -E_\theta\left[\frac{\partial^2}{\partial \theta_i \partial \theta_j} \log f(Y; \theta)\right]$$

## Theorem

Let  $Y = (Y_1, \dots, Y_n) \sim F(\cdot; \theta)$  Fisher regular dist

Let  $t(Y)$  be a (biased) estimator of  $\theta$ .

$$\text{Then } \text{Var}[t(Y)] \geq \frac{(\partial \mathbb{E}[t(Y)])^2}{I(\theta)} = \text{CRLB}(\theta)$$

If  $t(\cdot)$  is unbiased, then we get  $\text{Var}[t(Y)] \geq \frac{1}{I(\theta)}$

## MLE - Theory

## Asymptotic Normality

### Idea

The MLE is popular for multiple reasons, one of such being that MLE is asymptotically efficient:

in the limit, the MLE achieves the CRLB. Recall that the MLE and other point estimator are themselves random variables. Therefore, a low-variance estimator  $\hat{\theta}_n$  estimates the true parameter  $\theta_0$  more precisely.

**Theorem** Assuming a Fisher-regular distribution with parameter  $\theta_0$  from which an iid. sample  $Y_1, \dots, Y_n$  is drawn. ? So we assume  $G(\cdot) \in \mathcal{F}$ . Then the MLE is asymptotically normally distributed with

$$\hat{\theta} \xrightarrow{D} N(\theta_0, I(\theta_0)^{-1})$$

## Maximum Likelihood in Misspecified Models

### Idea

We derived the previous results under the assumptions that our model is specified correctly, i.e.  $G(\cdot) \in \mathcal{F}$ .

If we drop this assumption the identity  $E[s(\theta_0)] = \int \frac{\partial \ell(\theta_0)}{\partial \theta} dG(y)$  still holds, but the Fisher information

$$I(\theta_0) := -\int \frac{\partial^2 \ell(\theta_0)}{\partial \theta^2} dG(y)$$

is no longer equal to the variance of the score  $V(\theta_0) := \int s(\theta_0) \cdot s'(\theta_0) dG(y) = \text{Var}[s(\theta_0)]$ .

This changes the asymptotic normality behaviors of the MLE.

**Theorem** Assuming Fisher-regular and iid. sample  $Y_1, \dots, Y_n$  drawn from  $Y_i \sim G(\cdot)$ .

The MLE is asymptotically normally distributed with  $\hat{\theta} \xrightarrow{D} N(\theta_0, I''(\theta_0) \cdot V(\theta_0) \cdot I''(\theta_0))$

? Since  $G(\cdot)$  is unknown, neither  $I(\theta_0)$  nor  $V(\theta_0)$  can be calculated analytically. But empirical estimates are

$$\hat{I}(\theta_0) = \frac{1}{n} \sum_{i=1}^n \frac{\partial^2 \log(f(y_i; \hat{\theta}))}{\partial \theta^2}$$

$$\text{and } \hat{V}(\theta_0) = \frac{1}{n} \sum_{i=1}^n \left( \frac{\partial \log(f(y_i; \hat{\theta}))}{\partial \theta} \right)^2. \text{ If } G \in \mathcal{F}, \text{ then } V(\theta_0) = I(\theta_0).$$

**Parameter Transformation**

Let  $y = h(\theta)$  for some bijective transformation  $h$  with  $\hat{y} = h(\hat{\theta})$ .  $I_y(y) := I_\theta(h^{-1}(y))$ .

$$I_y(y) = E\left[-\frac{\partial^2 \log(f(y; \hat{\theta}))}{\partial y^2}\right] = \frac{\partial \theta}{\partial y} \cdot I_\theta(\theta) \cdot \frac{\partial \theta}{\partial y}$$

It follows  $\hat{y} \xrightarrow{D} N(y, \frac{\partial \theta}{\partial y} \cdot I_\theta(\theta) \cdot \frac{\partial \theta}{\partial y})$

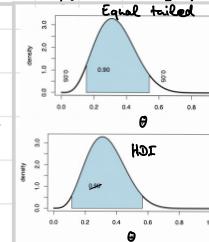
## Example

$$X_i \sim U(0, \theta), \text{ i.i.d. } t(X) = \max\{X_1, \dots, X_n\}$$

$$\begin{aligned} P(X_1 = x_1, \dots, X_n = x_n; \theta) &= \frac{1}{\theta^n} \cdot I_{(0, \infty)}(\max\{x_1, \dots, x_n\}) \cdot I_{(0, \infty)}(\min\{x_1, \dots, x_n\}) \\ &= g(t(X), \theta) = h(x_1, \dots, x_n) \end{aligned}$$

So the sample maximum is a sufficient statistic

for the population maximum.





## Nested and non-nested models

We look at two possible models which we denote as  $F_1$  and  $F_2$ . We prefer the model with the smaller Kullback-Leibler divergence,

If one model is a subset of the other, then we call the models nested, i.e. we choose model  $F_1$  over model  $F_2$  if

Assuming  $F_1$  to be the larger model this can be written as

## Model comparison idea

$$KL(F_1(\cdot) | F_2(\cdot; \hat{\theta}_2)) - KL(F_1(\cdot) | F_1(\cdot; \hat{\theta}_1)) < 0 \Leftrightarrow \int \log \frac{f_1(y; \hat{\theta}_1)}{f_2(y; \hat{\theta}_2)} dG(y) > 0$$

The quantity  $\sum_{i=1}^n \log \frac{f_1(y_i; \hat{\theta}_1)}{f_2(y_i; \hat{\theta}_2)}$  does not converge to the integral

because we use the data twice: ones for  $y_1, \dots, y_n$  and ones for estimating  $\hat{\theta}_1, \hat{\theta}_2$ .

Alternatively, we can write non-nested models formally as

$F_1 = \{F(\cdot; \theta) | \theta \in \Theta_1\}$  and  $F_2 = \{F(\cdot; \theta) | \theta \in \Theta_2\}$  with  $\Theta_2 \subsetneq \Theta_1$ . This would lead to overfitting, b.c. more complex models are preferred.

## Known Distribution Function

## Known Density Function (Rejection Sampling)

We want to draw random samples from a known distribution  $F(\cdot)$ . We want to draw random samples from an unknown distribution  $F$  by using its known density function  $f$ .

We need the inverse of  $F(\cdot)$ . We define  $F^{-1}(u) := \inf\{y | F(y) \geq u\}$ . We use rejection sampling to do so: Assume we know a distribution  $F^*$ , its inverse and its density  $f^*$ . Assume further that  $\exists a \in \mathbb{R}, \forall y \in T_F: f(y) \leq a \cdot f^*(y)$  and that  $a$  is known.

To sample from  $F(\cdot)$  we use the following property:

Let  $Y = F^{-1}(U)$  where  $U \sim U(0,1)$  then  $Y \sim F(\cdot)$ .

## Simulating Distr.

Then the simulation procedure is as follows:

1. Draw  $Y^*$  from  $F^*$  using the method described before.
2. Draw  $U$  from  $U(0,1)$ .
3. If  $U \leq \frac{f(Y^*)}{a \cdot f^*(Y^*)}$  accept  $Y^*$ . Otherwise go back to 1. ! acceptance probability is  $\frac{1}{a}$ .

We can show that the accepted values  $Y^*$  follow the distribution  $F$ .

## Metropolis (-Hastings) Algorithm

1. Select a starting value  $y^{(0)}$  and set  $y^{(t)} = y^{(0)}$

2. Based on  $y^{(t)}$  propose a new value  $y^*$  from a proposal distribution  $H(y^* | y^{(t)})$

where  $H(\cdot | y^{(t)})$  is a known distribution, with the same support (Trager) as  $F$ , from which we can sample. The corresponding density is  $h(y^* | y^{(t)})$ . In the Metropolis-Hastings algorithm the proposal distribution is symmetric - for example  $N(y^{(t)}, \sigma^2)$  with some fixed variance  $\sigma^2$ . In the Metropolis-Hastings

Algorithm the proposal distribution can be skewed.  $H(\cdot | y^{(t)})$  can depend on  $y^{(t)}$  but doesn't need to.

3. Compute the acceptance probability, which is defined as

$$\alpha(y^*, y^{(t)}) = \min \left\{ 1, \frac{f(y^*)}{f(y^{(t)})} \cdot \frac{h(y^{(t)} | y^*)}{h(y^* | y^{(t)})} \right\}.$$

If  $H$  is symmetric, then  $\frac{h(y^{(t)} | y^*)}{h(y^* | y^{(t)})} = 1$ .

$$4. \text{Draw } U^* \sim U(0,1) \text{ and define } y^{(t+1)} = \begin{cases} y^*, & \text{if } U^* \leq \alpha(y^*, y^{(t)}) \\ y^{(t)}, & \text{otherwise} \end{cases}$$

5. Go back to step 2 until  $t$  is large enough.

(6) Usually the first few values of  $y^{(t)}$ , i.e.  $y^{(1)}, \dots, y^{(500)}$  are discarded (Burn-in).

## Gibbs Sampling

Note that  $P_{\theta|y_1, \dots, y_n}(\theta | y_1, \dots, y_n) \propto \prod_{i=1}^n f(y_i; \theta) \cdot p(\theta)$ .

1. Let  $\Theta = (\theta_1, \dots, \theta_p) \in \Theta$ . Let  $\theta_j \in \Theta$  and set  $(\theta_1^{(0)}, \dots, \theta_p^{(0)}) = \theta^{(0)} = \theta_0$ .

2. Set  $\theta_j^{(t+1)} = \theta_j^{(t)}$  for  $j=1, \dots, p$  and execute the following steps:

(i) Set  $\theta^* = \theta^{(t)}$

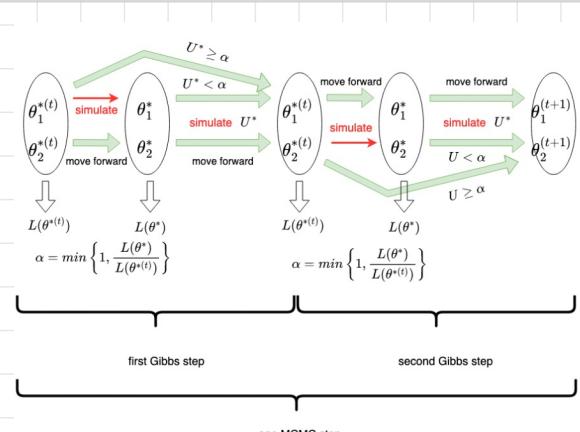
(ii) Replace  $\theta_j^*$  with a sample from the proposal  $q(\cdot; \theta_j^{(t)})$

(iii) Sample  $U^* \sim U(0,1)$

(iv) If  $U^* \leq \alpha(\theta_j^{(t)}, \theta_j^{(t+1)}) = \min \left\{ 1, \frac{L(\theta^{(t+1)})}{L(\theta^{(t)})} \cdot \frac{q(\theta_j^{(t)} | \theta_j^{(t+1)})}{q(\theta_j^{(t+1)} | \theta_j^{(t)})} \right\}$ , then

set  $\theta_j^{(t+1)} = \theta_j^{(t)}$ . Otherwise leave  $\theta_j^{(t+1)}$  unchanged.

3. Set  $\theta^{(t+1)} = \theta^{(t)}$  and repeat step 2 until a stationary dist. is reached.



## Akaike Information Criterion (AIC)

## AIC/BIC Theory

The AIC for model class  $\mathcal{F} = \{F(\cdot; \theta) | \theta \in \Theta\}$  with  $p = \dim(\Theta)$

$$\text{is defined as } AIC(\mathcal{F}) := -2 \sum_{i=1}^n \log(f(y_i; \hat{\theta})) + 2p.$$

The AIC is an estimate for the quantity  $2 \cdot E[KL(F(\cdot), F(\cdot; \hat{\theta}))] - 2 \cdot \log(f(y_i; \hat{\theta}))$ .

Therefore  $AIC(\mathcal{F}_1) - AIC(\mathcal{F}_2)$  is an estimate for  $-2 \cdot E[\log \frac{f_1(y; \hat{\theta}_1)}{f_2(y; \hat{\theta}_2)} dG(y)]$ .

Hence, if  $AIC(\mathcal{F}_1) < AIC(\mathcal{F}_2)$ , we prefer  $\mathcal{F}_1$  and otherwise  $\mathcal{F}_2$ .

## Bayesian Information Criterion (BIC)

The idea is the same as with AIC but the BIC of a model

$$\mathcal{F} \text{ is defined as } BIC(\mathcal{F}) := -2 \cdot \sum_{i=1}^n \log(f(y_i; \hat{\theta})) + \log(n) \cdot p.$$

The BIC favors models with lower complexity, i.e. fewer parameters.

## Corrected AIC (AICc)

In the case of small sample statistics the AIC needs to be modified and one should use the corrected AIC defined as  $AICc := -2 \cdot \ell(\hat{\theta}) + 2p \left( \frac{n}{n-p-1} \right)$ .

## Importance Sampling

Same setup as in rejection

sampling but we are interested in

$$E_{\tilde{F}}[h(Y)] \text{ for some function } h(\cdot).$$

$$\text{We use } E_{\tilde{F}}[h(Y)] = \int_R h(y) f(y) dy = \int_R h(y) \frac{f(y)}{f^*(y)} f^*(y) dy = E_{F^*}[h(y) \frac{f(y)}{f^*(y)}].$$

Using now i.i.d. samples  $y_1, \dots, y_n \sim F^*(\cdot)$

$$\text{we obtain } \frac{1}{N} \sum_{i=1}^N h(y_i) \frac{f(y_i)}{f^*(y_i)} \xrightarrow{N \rightarrow \infty} E_{\tilde{F}}[h(Y)].$$

! The closer  $f$  is to  $f^*$  the smaller the variance of our approximation.

## AIC, BIC and Hypothesis Testing

To connect the ideas of model selection with hypothesis testing we assume two nested models  $\mathcal{F}_0$  and  $\mathcal{F}_1$  with  $\mathcal{F}_0 \subset \mathcal{F}_1$ . Assume  $\theta_0 = |\Theta_0|$  and  $\theta_1 = |\Theta_1|$  and  $p = \theta_1 - \theta_0 > 0$ .

Let  $\hat{\theta}_0$  and  $\hat{\theta}_1$  be the corresponding MLEs. It can be shown that if model  $\mathcal{F}_0$  holds, then

$$2(\ell(\hat{\theta}_0) - \ell(\hat{\theta}_1)) \sim \chi_p^2.$$

Formulating the hypothesis  $H_0: \mathcal{F} = \mathcal{F}_0$  and alternative  $H_1: \mathcal{F} = \mathcal{F}_1$ , we obtain the decision rule " $H_1$ "  $\Leftrightarrow 2(\ell(\hat{\theta}_0) - \ell(\hat{\theta}_1)) > \chi_{p, 1-\alpha}^2$ . On the other hand, the decision rule of the AIC can be written as " $H_1$ "  $\Leftrightarrow 2(\ell(\hat{\theta}_0) - \ell(\hat{\theta}_1)) > 2p$  and for the BIC we get " $H_1$ "  $\Leftrightarrow 2(\ell(\hat{\theta}_0) - \ell(\hat{\theta}_1)) > \log(n)p$

## Prior and Posterior

Let  $\mathcal{F} = \{f(\cdot; \theta) | \theta \in \Theta\}$  be our probability model for the i.i.d. data  $y_1, \dots, y_n$ .

For  $\theta$  we formulate our (missing) knowledge as prior distribution  $\theta \sim p(\cdot; \gamma)$  where parameter  $\gamma \in \Gamma$  is called hyper-parameter. We define the prior-structure as the model class  $\mathcal{P} = \{p(\cdot; \theta) | \theta \in \Theta\}$ .

The posterior distribution  $p_{\theta|y_1, \dots, y_n}(\theta; y_1, \dots, y_n) = \frac{\prod_{i=1}^n f(y_i; \theta) p(\theta; \gamma)}{\int_{\Theta} \prod_{i=1}^n f(y_i; \theta) p(\theta; \gamma) d\theta}$  with  $f(y; \gamma) = \int_{\Theta} f(y; \theta) p(\theta; \gamma) d\theta$ .

## Bayes Inference

## Conjugate prior distribution

For a given family of distributions  $\mathcal{F}$ , we call  $\mathcal{P}$  the set of conjugate prior distributions, if it exists and if  $p_{\theta|y_1, \dots, y_n}(\theta; y_1, \dots, y_n) \in \mathcal{P}$ . That is the posterior is from the same family of distributions as the prior distribution.

## Example

Lets assume  $Y \sim Po(1)$  and  $\lambda \sim \Gamma(\alpha, \beta)$ , i.e.  $p(\lambda) \propto \lambda^{\alpha-1} \exp(-\lambda \beta)$ .

$$\text{Then } p_{\lambda|y_1, \dots, y_n}(\lambda; y_1, \dots, y_n) \propto \left( \prod_{i=1}^n \lambda^{y_i} \exp(-\lambda) \right) \cdot \lambda^{\alpha-1} \exp(-\lambda \beta) = \lambda^{\frac{1}{2} \sum y_i + \alpha - 1} \exp(-\lambda(\alpha + \beta)).$$

Consequently  $\lambda | y_1, \dots, y_n \sim \Gamma(\frac{1}{2} \sum y_i + \alpha, \alpha + \beta)$ . So given the set of Poisson distribution  $\mathcal{F}$  with parameter  $\lambda$ ,

the set of gamma distributions  $\mathcal{P}$  is a set of conjugate prior distributions.

## Bernstein-von Mises Theorem

For increasing sample size  $n$  and appropriately chosen prior we find  $\theta \sim N(\hat{\theta}, I^{-1}(\hat{\theta}))$

Parameter	$\mathcal{F}$	$\mathcal{P}$
$\pi$	Binomial distribution	Beta distribution
$\lambda$	Poisson distribution	Gamma distribution
$\mu$	Normal distribution	Normal distribution
$\lambda$	Exponential distribution	Gamma distribution

Table 10.1 Examples for conjugate distributions

## Approximate Bayes Computation (ABC)

### Concept

To simulate from the posterior  $\theta | y_1, \dots, y_n$  by using MCMC

we need to know/compute the likelihood  $f(y_i; \theta)$ .

If the likelihood is unknown/uncomputable, then we can use the

ABC -method to at least approximately simulate from the posterior.

### Exact Algorithm

1. Generate  $\theta^*$  from some prior dist.  $p_0(\theta)$ .
2. Generate  $y^*$  from the assumed dist. i.e.  $f_{y^*}(\cdot; \theta^*)$ .
3. If  $y^* = y$ , with  $y$  being the available data,
4. Then accept  $\theta^*$ . Otherwise go back to step 1.
5. Repeat 1-3, until a sufficient amount of simulated values is reached.

### Approx. Algorithm

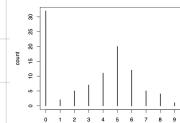
1. Generate  $\theta^*$  from some prior dist.  $p_0(\theta)$ .
2. Generate  $y^*$  from the assumed dist. i.e.  $f_{y^*}(\cdot; \theta^*)$ .
3. If  $d(y^*, y) \leq \epsilon$ , with  $y$  being the available data, then accept  $\theta^*$ . Otherwise go back to step 1.  
 $\text{Loc}(\cdot, \cdot)$  is a distance measure and  $\epsilon$  needs to be chosen appropriately.
4. Repeat 1-3, until a sufficient amount of simulated values is reached.

### Example - Zero inflated data

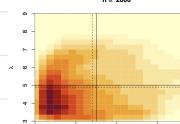
Consider the same setup as above. We define  $d(y^*, y) := \sum_{k=0}^N (y_k^* - x_k)^2$   
with  $N = 10$ ,  $\bar{x}_k := \frac{1}{n} \sum_{i=1}^n \mathbb{1}(y_i = k)$  and  $\bar{x}_k^* := \frac{1}{n} \sum_{i=1}^n \mathbb{1}(y_i^* = k)$  and  $\epsilon = 0.05$ .

Now we proceed with the algorithm and simulate  $\theta^*, \lambda^*$ .

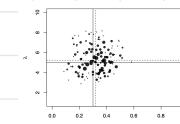
Actual data  $y_1, \dots, y_n$



Posteriori with MCMC



Posteriori with ABC



Point size  $\propto d(y^*, y)$

## Variational Bayes Reasoning

### Concept

In this method we have the same setting and goal as in MCMC, but we trade the accuracy of MCMC (exact posterior) for speed/computation time. Instead of simulating exactly from the posterior we simulate from a more simple distribution  $q(\cdot)$  that is approximately the posterior  $p(\theta|y)$ , e.g.  $q(\cdot) = N(\hat{\theta}_p, I_p(\hat{\theta}_p))$  (see Thm. below).

In general, let  $q(\theta, y) \in Q = \{q(\cdot, y) | y \in \Gamma\}$  be some easy-to-handle dist. and we are interested in  $q(\cdot; \hat{\theta}) = \arg \min_{q \in Q} \{KL(q(\cdot, \cdot), p(\cdot|y))\}$ . This results in a less accurate estimate compared to MCMC but is less computationally expensive than MCMC.

### Solving the Optimization Problem

$$KL(q(\cdot, y), p(\cdot|y)) = \int \log \left( \frac{q(\theta, y)}{p(\theta|y)} \right) q(\theta, y) d\theta = - \underbrace{\int \log \left( \frac{f(y|\theta) p(\theta)}{q(\theta, y)} \right) q(\theta, y) d\theta}_{\text{unknown}} + \log(f(y))$$

To minimize  $KL(q(\cdot, y), p(\cdot|y))$  we need to maximize  $LB(q(\cdot, y))$ .

$$\frac{\partial LB(q(\cdot, y))}{\partial \theta} = \int \log(f(y|\theta)) \frac{\partial g(\theta, y)}{\partial \theta} d\theta - \int \log(q(\theta, y)) \frac{\partial g(\theta, y)}{\partial \theta} d\theta - \int \frac{\partial \log(q(\theta, y))}{\partial \theta} q(\theta, y) d\theta = \mathbb{E}_{q(\cdot, y)} \left[ \frac{\partial \log(q(\theta, y))}{\partial \theta} \log \left( \frac{f(y|\theta) p(\theta)}{q(\theta, y)} \right) \right] = 0 \quad (\text{see Score Function})$$

! Because of the unknown component we can't quantify how close our approximation actually is.

### Approximating Solution of the Optimization Problem

$$\text{We want to solve } \frac{\partial LB(q(\cdot, y))}{\partial \theta} = \mathbb{E}_{q(\cdot, y)} \left[ \frac{\partial \log(q(\theta, y))}{\partial \theta} \log \left( \frac{f(y|\theta) p(\theta)}{q(\theta, y)} \right) \right] = 0.$$

We can replace the expected value by simulated values, i.e. for given  $y^{(t)}$  we

draw  $\theta^{(k)} \sim q(\cdot; y^{(t)})$  for  $k = 1, \dots, K$  and compute

$$\frac{\partial LB(q(\cdot, y^{(t)}))}{\partial \theta} = \frac{1}{K} \sum_{k=1}^K \frac{\partial \log(q(\theta^{(k)}, y^{(t)}))}{\partial \theta} \cdot \log \left( \frac{f(y^{(t)}|\theta^{(k)}) p(\theta^{(k)})}{q(\theta^{(k)}, y^{(t)})} \right)$$

We can now use  $y^{(t+1)} = y^{(t)} + \frac{\partial LB(q(\cdot, y^{(t)}))}{\partial \theta}$  to solve the problem numerically.  
↑ learning rate

### Penalized version of Bernstein Thm.

For sufficiently large  $n$  we have  $\theta|y \approx N(\hat{\theta}_p, I_p(\hat{\theta}_p))$

$$\text{where } \hat{\theta}_p := \arg \max_{\theta \in \Theta} l_p(\theta) = \arg \max_{\theta \in \Theta} \sum_{i=1}^n \log(f(y_i|\theta)) + \log(p(\theta))$$

where the prior  $p$  might depend on further hyperparameters and  $I_p(\theta) := I(\theta) + \frac{\partial^2 p(\theta)}{\partial \theta^2}$ .

### Log-Linear Model (Multiple observation)

Let  $Y^{(1)}, \dots, Y^{(n)}$  be independent observations, i.e.  $Y^{(i)} = (Y_1^{(i)}, \dots, Y_p^{(i)})$ .

$$P(Y^{(1)} = y^{(1)}, \dots, Y^{(n)} = y^{(n)}) \propto \prod_{i=1}^n \exp \left( \theta_0 + \sum_{j=1}^{K_1} \sum_{k_j=1}^{K_{j+1}} s_{1, k_1, \dots, j, k_j} (y^{(i)}) \cdot \theta_{k_1, \dots, k_j} + \dots + \sum_{j=2}^{p-1} \sum_{k_j=1}^{K_{j+1}} s_{j, k_1, \dots, k_{j-1}, k_j} (y^{(i)}) \cdot \theta_{k_1, \dots, k_{j-1}, k_j} + \dots + \sum_{k_{p-1}=1}^{K_p} s_{p, k_{p-1}} (y^{(i)}) \cdot \theta_{k_{p-1}} \right)$$

$$= \exp \left( \theta_0 + \sum_{j=1}^{p-1} \sum_{k_j=1}^{K_{j+1}} n_{j, k_1, \dots, k_j} \cdot \theta_{k_1, \dots, k_j} + \sum_{j=2}^{p-1} \sum_{k_j=1}^{K_{j+1}} s_{j, k_1, \dots, k_{j-1}, k_j} (y^{(i)}) \cdot \theta_{k_1, \dots, k_{j-1}, k_j} + \dots + \sum_{k_{p-1}=1}^{K_p} s_{p, k_{p-1}} (y^{(i)}) \cdot \theta_{k_{p-1}} \right)$$

zero variables      sums over one variable      sums over two-wise pairs of variables      sum over all p variables

We call the terms  $\theta_{\dots}$  with more than one variable involved interaction terms, e.g.  $\theta_{k_1, \dots, k_p}$  is the pairwise interaction between  $Y_1$  and  $Y_p$  for the outcome  $k_1, \dots, k_p$ .

! The interaction term is zero iff the involved variables are independent.

## Multivariate Data

### 3-dimensional Log-Linear model and model hierarchy

We look at a Log-linear model for count-data with  $p=3$  and  $K_1 = K_2 = K_3 = 2$ . Let  $Y = (Y_1, \dots, Y_p)$  be a  $p$ -dimensional random variable with  $f_{123}: p(y_1, \dots, y_p)$  as corresponding density or prob. function. We define  $A, B, C \subseteq \{1, \dots, p\}$  s.t.

### Conditional Independence

$A \perp B = A \cap B = \emptyset$  and  $A \perp \emptyset$  and  $B \perp \emptyset$ .

We say  $A$  and  $B$  are conditionally independent given  $C$ , denoted as  $(A \perp B) | C$

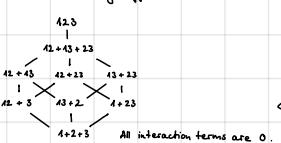
iff  $f_{(A \perp B) | C}(y_A, y_B | y_C) = f_{AC}(y_A | y_C) \cdot f_{BC}(y_B | y_C)$ . This holds iff

$$f_{ABC}(y_A, y_B, y_C) = h_{AC}(y_A | y_C) \cdot h_{BC}(y_B | y_C) \text{ for some function } h_{AC}, h_{BC}.$$

If we impose  $\theta_{mn} = 0$ , then we can denote the model as  $1+2+1+3+2+3$  or  $42+43+23$ .

If we further restrict  $\theta_{mn} = 0$ , then we can denote the model as  $1+3+2+3$  or  $13+23$ .

If we continue on setting different interaction terms to 0, then this leads to some model-hierarchy:



Let  $A, B, C$  be defined as before. We can visualise  $(A \perp B) | C$  as the graph

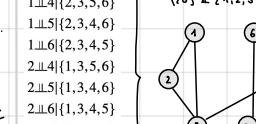
The missing edge between  $A, B$  and the path that goes through  $C$  translates to  $(A \perp B) | C$ .

If instead of sets of variables we want to visualise single variables, we construct the graph as follows: Let  $G = (E, V)$  be a graph, i.e.

$E = \{1, \dots, p\}$  and  $V = E \times E$ . Let  $i, j \in E$ .  $i \perp j | E \Leftrightarrow (i, j) \in V$

### Example

From the graph we can see  
 $\{1, 2\} \perp \{4, 5, 6\} | \{3\}$   
 $\{5, 6\} \perp \{1, 2, 3, 4, 5\} | \{3\}$



### Multivariate Normal Dist.

$$\text{Let } Y \sim \mathcal{N}_p(\mu, \Sigma)$$

$$f(y; \mu, \Sigma) = \frac{\exp \left( -\frac{1}{2} (y - \mu)^T \Sigma^{-1} (y - \mu) \right)}{\sqrt{(2\pi)^p \det(\Sigma)}}$$

Using  $\Sigma = \Sigma^{-1}$  we can rewrite  $f$  as

The inverse of the covariance matrix

$$f(y; \mu, \Sigma) \propto \exp \left( -\frac{1}{2} \sum_{i=1}^p \sum_{j=1}^p y_i^2 + \frac{1}{2} \sum_{j=1}^p \sum_{i=1}^p y_i y_j \right)$$

If  $\Sigma_{ij} = 0$  and  $C = \{1, \dots, p\} \setminus \{j\}$ , then

$$f(y; \mu, \Sigma) \propto h_{C \cup \{j\}}(y_C, y_j) \cdot h_{C \cup \{j\}}(y_C, y_j)$$

Hence,  $\Sigma_{je} = 0 \Leftrightarrow (Y_e \perp Y_j) | Y_C$

### Conditional Independence in Log-Linear Model

Let  $Y = (Y_1, \dots, Y_p)$  be a discrete random variable with  $Y_i \in \{1, \dots, K_i\}$   $i = 1, \dots, p$ .

$$(Y_1 \perp Y_2) | Y_3 \text{ or } (Y_1 \perp Y_2) | Y_3 \Leftrightarrow \text{Var}_{Y_1, Y_2, Y_3} \text{V} \in \{1, \dots, K_1, K_2\} \cdot \frac{P(Y_1 = k_1 | Y_3 = k_3) \cdot P(Y_2 = k_2 | Y_3 = k_3)}{P(Y_1 = k_1, Y_2 = k_2 | Y_3 = k_3)} = \frac{h_{123}(k_1, k_2, k_3)}{h_{13}(k_3)}$$

$$\Leftrightarrow \text{Var}_{Y_1, Y_2, Y_3} \text{V} \in \{1, \dots, K_1, K_2\} \cdot P(Y_1 = k_1 | Y_3 = k_3) \cdot P(Y_2 = k_2 | Y_3 = k_3) \quad (\text{given } Y_3 \text{ we don't get additional information from } Y_2)$$

In our log-linear model, this holds iff all parameters that involve components  $Y_1, Y_2$  are 0.

i.e.  $(Y_1 \perp Y_2) | Y_3 \Leftrightarrow \text{Var}_{Y_1, Y_2, Y_3} \text{V} \in \{1, \dots, K_1, K_2\} \cdot \theta_{k_1 k_2} = 0$ .

Using the model notation from before, this results in the model  $13+23$ .

## Rare Events

### Setup

Let  $Y \sim B(n, \pi)$  and we want to estimate  $\pi$  based on the sample  $y_{1:n}, y_{2:n}$ .

But if  $\pi = 0$  or in general  $n\pi = 0$  then our maximum likelihood estimation method fails and the usual confidence interval is meaningless.

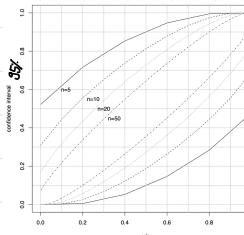
Therefore we need another way to construct a confidence interval for rare events.

### Clopper and Pearson Exact Confidence Interval

Let  $Y \sim B(n, \pi)$ . For  $Y=y$  the exact confidence interval is  $(c_L, c_U)$  with  $c_L := \inf \{ \pi \mid P(Y \leq y, \pi) \geq \frac{\alpha}{2} \}$ ,  $c_U := \sup \{ \pi \mid P(Y \geq y, \pi) \geq \frac{\alpha}{2} \}$ .

This can be computed by using the Beta-distribution:

$$c_L = Be\left(\frac{\alpha}{2}, y, n-y+1\right), \quad c_U = Be\left(1-\frac{\alpha}{2}, y+1, n-y\right) \quad (\text{in R: pbeta})$$



## Generalized Extreme Value Distribution

The GEV distribution has three parameters.

location parameter  $\mu$ , scale parameter  $\sigma$ , shape parameter  $\gamma$ .

The distribution function is given by

$$H(z) = \begin{cases} \exp(-\exp(-z)) & \text{for } \gamma = 0 \\ \exp(-(1+z/\sigma)^\gamma) & \text{for } \gamma \neq 0 \end{cases}$$

$$1. \text{ Gumbel distribution } (\gamma=0): F(x; \mu, \sigma, 0) = \exp(-\exp(-\frac{x-\mu}{\sigma}))$$

$$2. \text{ Fréchet distribution } (\gamma > 0): F(x; \mu, \sigma, \gamma) = \exp(-(1 + \frac{x-\mu}{\sigma})^{-\gamma}) \cdot \mathbb{1}_{(0, \infty)}(1 + \frac{x-\mu}{\sigma})$$

$$3. \text{ (Reversed) Weibull distr. } (\gamma < 0): F(x; \mu, \sigma, \gamma) = \begin{cases} \exp(-(1 + \frac{x-\mu}{\sigma})^{-\gamma}) & \text{if } 1 + \frac{x-\mu}{\sigma} > 0 \\ 1 & \text{otherwise} \end{cases}$$

Assume  $\mu_y(t) = 0$  and  $\gamma_y(t, t') = \gamma_y(|t-t'|) = \gamma_y(h) = \text{Cov}(y_t, y_{t+h})$  with  $h \in \{0, 1, \dots, T\}$ .

$\gamma_y(h)$  is called autocovariance function.

$$\text{The variance matrix is defined as } \Gamma = \begin{bmatrix} \gamma(0) & \gamma(1) & \dots & \gamma(T) \\ \gamma(1) & \gamma(0) & \dots & \vdots \\ \vdots & \vdots & \ddots & \gamma(0) \\ \gamma(T) & \dots & \gamma(0) & \gamma(0) \end{bmatrix}.$$

B.c. of  $\gamma(0) = V(y_t)$  we know that  $\rho(1) = \frac{\gamma(1)}{\gamma(0)}$ . We call  $\rho(h)$  the autocorrelation function

$$\text{The correlation matrix is defined as } R = \begin{bmatrix} \rho(0) & \rho(1) & \dots & \rho(T) \\ \rho(1) & \rho(0) & \dots & \rho(1) \\ \vdots & \vdots & \ddots & \vdots \\ \rho(T) & \dots & \rho(1) & \rho(0) \end{bmatrix}.$$

We define white noise as the time series  $y_t \sim N(0, \sigma^2)$  with  $\gamma_y(h) = 0$  for  $h \neq 0$

### MA(1)

$$y_t = z_t + \theta z_{t-1}. \text{ Then } \gamma(0) = \sigma^2(1+\theta^2)$$

$$\gamma(h) = \theta^h z_{t-h} z_t = \theta^h \gamma(0) \text{ for } h > 1 \text{ and hence } \rho(h) = \begin{cases} 1, & h=0 \\ 0, & h>1 \end{cases}$$

$$y_t = \phi z_t + \varepsilon_t. \text{ Then } \gamma(0) = \phi^2(1+\sigma^2)$$

$$\gamma(h) = \phi^h \varepsilon_{t-h} \varepsilon_t = \phi^h \gamma(0) \text{ for } h > 1$$

$$\text{This results in } V[y_t] = \sum_j \gamma(j) \sigma^2 = \frac{\sigma^2}{1-\phi^2}$$

### Datengleichheit

$$\text{Let } \mu = E[Y] \text{ and } \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i = \frac{\sum_{i=1}^n R_i y_i^{(i)}}{\sum_{i=1}^n R_i^{(i)}}.$$

$$\mu - \bar{y} = \underbrace{\text{Qualität}}_{\text{Corr}(R, Y)} \cdot \sqrt{\frac{\text{Var}(Y)}{\text{Variabilität}}} \cdot \sqrt{\frac{n-n}{n}}$$

### Auto Regressive Process AR(p)

$$y_t = \phi_1 y_{t-1} + \phi_2 y_{t-2} + \dots + \phi_p y_{t-p} + z_t \text{ with } z_t \sim N(0, \sigma^2) \text{ i.i.d.}$$

We can rewrite this using the backshift operator  $B$ , i.e.  $B^k y_t = y_{t-k}$ ,  $k \in \mathbb{Z}$ . Hence  $\rho(1) = \frac{\phi_1 \gamma(0)}{\gamma(0)} = \phi$ . In general  $\rho(h) = \phi^h$ . For  $|\phi| < 1$  we get

$$\text{The AR}(p) \text{ process can then be written as } (1 - \sum_{j=1}^p \phi_j B^j) y_t = z_t$$

### AR(1)

$$y_t = \phi y_{t-1} + z_t. \text{ Then } \text{Cov}(y_t, y_{t-h}) = \phi \text{Cov}(y_{t-1}, y_{t-h}) = \phi \gamma(h)$$

$$\text{We have } \gamma(0) = V(y_t) = \phi^2 \gamma(0) + \sigma^2 \text{ and } \gamma(h) = \phi^h \gamma(0).$$

$$\text{The variance matrix is defined as } \Gamma = \begin{bmatrix} \gamma(0) & \gamma(1) & \dots & \gamma(T) \\ \gamma(1) & \gamma(0) & \dots & \vdots \\ \vdots & \vdots & \ddots & \gamma(0) \\ \gamma(T) & \dots & \gamma(0) & \gamma(0) \end{bmatrix}.$$

$$\text{Note that } y_t = \phi y_{t-1} + z_t = \phi^2 y_{t-2} + \phi z_{t-1} + \phi^2 y_{t-2} + z_t = \dots = \sum_{j=0}^{T-1} \phi^j z_{t-j}$$

$$\text{This results in } V[y_t] = \sum_j \gamma(j) \sigma^2 = \frac{\sigma^2}{1-\phi^2}$$

### Causal process

An ARMA(p,q) is causal if it can be rewritten to

$$y_t = \sum_{j=0}^{\infty} \psi_j B^j z_{t-j} = \Psi(B) z_t$$

$$\text{with } \sum_j |\psi_j| < \infty$$

!

### ARMA(p,q)

$$y_t = \underbrace{\phi_1 y_{t-1} + \phi_2 y_{t-2} + \dots + \phi_p y_{t-p} + z_t}_{\text{AR}(p)} + \underbrace{\theta_1 z_{t-1} + \theta_2 z_{t-2} + \dots + \theta_q z_{t-q}}_{\text{MA}(q)}$$

$$\text{We define } \Phi(B) := -\sum_{j=0}^p \phi_j B^j \text{ and } \Theta(B) := \sum_{j=0}^q \theta_j B^j \text{ with } \phi_0 = 1, \theta_0 = 1$$

$$\text{Then we can rewrite the ARMA}(p,q) \text{ as } \Phi(B) y_t = \Theta(B) z_t$$

Representation is not unique with that notation. Can be solved if we find  $\Phi^{-1}(B)$  s.t.  $y_t = \Phi^{-1}(B) \cdot \Theta(B) z_t$ .

### Example

$$\text{ARMA}(1,1) \text{ process: } \Phi(B) = 1 - \phi B, \quad \Theta(B) = 1 + \theta B$$

$$\text{then } \Phi^{-1}(B) = \sum_{j=0}^{\infty} \phi^j B^j \text{ and } y_t = \sum_{j=0}^{\infty} \psi_j B^j z_{t-j}$$

$$\text{with } \psi_0 = 1 \text{ and } \psi_j = \phi^{j-1} (\phi + \theta)$$

Hence, all AR components can be transferred into MA components.

PACF for p

ACF for q

## EM - Algorithmus

1. Expectation step: Impose missing values based on  $\theta_{(t)}$

$$\text{Compute } Q(\theta, \theta_{(t)}) = \sum_{i=1}^n \int l_i(\theta) f(y_{1:n} | y_{(t)}, \theta_{(t)}) dy_{(t)}$$

with  $l_i(\theta) = \log(f(y_i; \theta))$  and  $y_{(t)} = (y_1, y_{(t)})$

! Easy for exponential family.

2. Maximisation step: Calculate the new MLE using the imputed values.

Maximise  $Q(\theta, \theta_{(t)})$  w.r.t.  $\theta$ , i.e. solve

$$s(\theta_{(t+1)}; \theta_{(t)}) = 0 \text{ with } s(\theta; \theta_{(t)}) = \frac{\partial Q(\theta, \theta_{(t)})}{\partial \theta}$$

3. Return to step 1 until convergence