

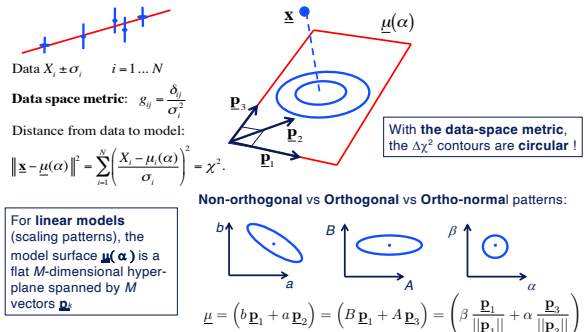
ADA 11 - 9am Thu 06 Oct 2022

Orthogonal Patterns (= orthogonal vectors using the data space metric)
e.g. Gram-Schmidt Orthogonalisation

Occam's Razor (model selection)
Information Criteria (AIC,BIC)

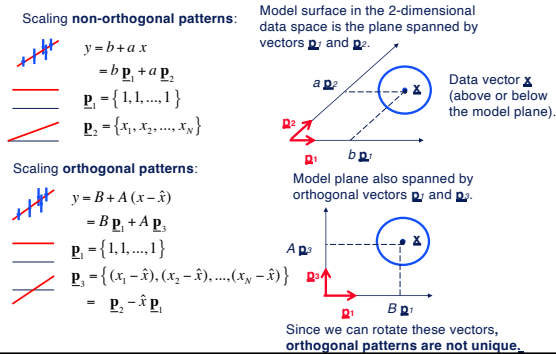
209

Review: Data Space Metric



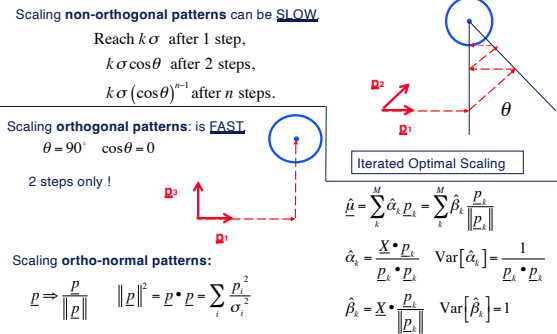
210

Scaling Orthogonal Patterns



211

Scaling Orthogonal Patterns



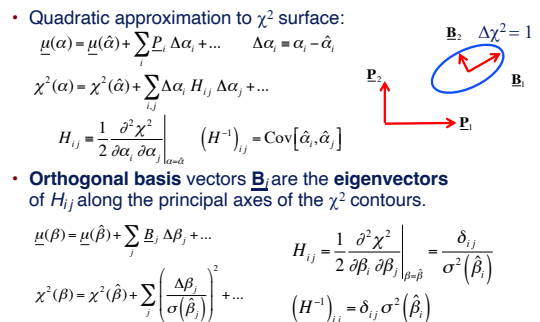
212

How to construct Orthogonal Patterns

- 1. Diagonalise Hessian Matrix
- 2. Gram-Schmidt Process
- 3. Differences between successive χ^2 fits

213

1. Diagonalise Hessian Matrix



214

Hessian Matrix for Non-Linear Models

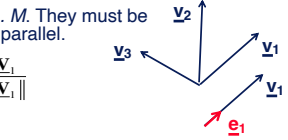
Model and derivatives: $\mu(\alpha), \frac{\partial \mu}{\partial \alpha_k} = \mathbf{P}_k, \frac{\partial^2 \mu}{\partial \alpha_k \partial \alpha_j} = \mathbf{C}_{kj}$
 M Gradient vectors: \mathbf{P}_k $M(M+1)/2$ Curvature vectors: \mathbf{C}_{jk}
 Badness-of-fit: $\chi^2 = \sum_{i=1}^N \left(\frac{y_i - \mu_i}{\sigma_i} \right)^2 = \|\mathbf{y} - \boldsymbol{\mu}\|^2$
 $\frac{\partial \chi^2}{\partial \alpha_k} = -2 \sum_{i=1}^N \frac{y_i - \mu_i}{\sigma_i^2} \frac{\partial \mu_i}{\partial \alpha_k} = -2(\mathbf{y} - \boldsymbol{\mu}) \cdot \mathbf{P}_k$
 Hessian Matrix: $H_{jk} = \frac{1}{2} \frac{\partial^2 \chi^2}{\partial \alpha_j \partial \alpha_k} = \mathbf{P}_j \cdot \mathbf{P}_k - (\mathbf{y} - \boldsymbol{\mu}) \cdot \mathbf{C}_{jk}$
 Best fit Parameters: $0 = \frac{\partial \chi^2}{\partial \alpha_k} = -2(\mathbf{y} - \boldsymbol{\mu}) \cdot \mathbf{P}_k \Rightarrow \boldsymbol{\mu}(\hat{\alpha}) \cdot \mathbf{P}_k = \mathbf{y} \cdot \mathbf{P}_k$
 Parameter Error Bars: $\text{Cov}[\hat{\alpha}_j, \hat{\alpha}_k] = (\mathbf{H}^{-1})_{jk}$
 Linear Model: $\mathbf{C}_{jk} = 0$ $H_{jk} = \mathbf{P}_j \cdot \mathbf{P}_k$

215

2. Gram-Schmidt Orthogonalization

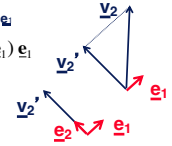
The **Gram-Schmidt process**:

- 1. Start with M vectors $\mathbf{V}_j, j = 1 \dots M$. They must be independent, i.e. no two of them parallel.
- 2. Normalize vector \mathbf{V}_1 : $\mathbf{e}_1 = \frac{\mathbf{V}_1}{\|\mathbf{V}_1\|}$



- 3. Make \mathbf{V}_2' perpendicular to \mathbf{e}_1 :
 - i.e. subtract component of \mathbf{V}_2 in direction of \mathbf{e}_1

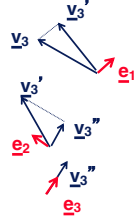
$$\mathbf{V}_2' = \mathbf{V}_2 - (\mathbf{V}_2 \cdot \mathbf{e}_1) \mathbf{e}_1$$



- 4. Normalize \mathbf{V}_2' : $\mathbf{e}_2 = \frac{\mathbf{V}_2'}{\|\mathbf{V}_2'\|}$

2. Gram-Schmidt Orthogonalization

- 5. Make \mathbf{V}_3' perpendicular to \mathbf{e}_1 :
 $\mathbf{V}_3' = \mathbf{V}_3 - (\mathbf{V}_3 \cdot \mathbf{e}_1) \mathbf{e}_1$
 - 6. Make \mathbf{V}_3'' perpendicular to \mathbf{e}_2 :
 $\mathbf{V}_3'' = \mathbf{V}_3' - (\mathbf{V}_3' \cdot \mathbf{e}_2) \mathbf{e}_2$
 - Note: \mathbf{V}_3'' is perpendicular to \mathbf{e}_1 AND \mathbf{e}_2 .
 - 7. Normalize \mathbf{V}_3'' : $\mathbf{e}_3 = \frac{\mathbf{V}_3''}{\|\mathbf{V}_3''\|}$
- ... Make \mathbf{V}_4 perpendicular to $\mathbf{e}_1, \mathbf{e}_2, \mathbf{e}_3$ and normalise to get \mathbf{e}_4 .
 Repeat up to \mathbf{V}_M to get **complete ortho-normal basis** $\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_M$.
 Easy to code ! (Try it !)

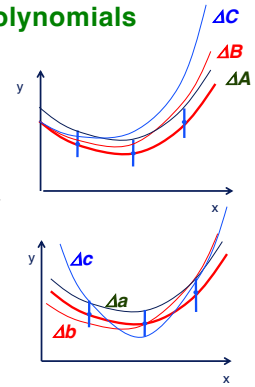


217

Orthogonal Polynomials

non-orthogonal polynomial:

$$y = A + Bx + Cx^2$$



orthogonal polynomials:

$$P_i \cdot P_j = \sum_{k=1}^N \frac{P_i(x_k) P_j(x_k)}{\sigma_k^2} = \frac{\delta_{ij}}{\text{Var}[\alpha_i]}$$

$$y = a P_0(x) + b P_1(x) + c P_2(x)$$

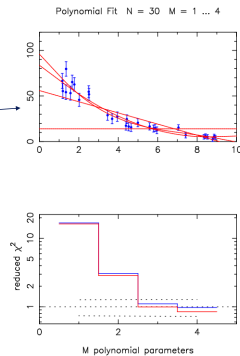


Note: every dataset has its own $1/\sigma^2$ weights, defines its own orthogonal polynomials.

218

3. Differences between successive χ^2 fits

- Fit: $A + Bx + Cx^2$
 - A, B, C are not independent
 - $1, x, x^2$ are not orthogonal
- If $P_k(x)$ is a polynomial of degree k fitted to the data, then $P_k(x) - P_{k-1}(x)$ are **orthogonal**:
- $a P_0(x) + b [P_1(x) - P_0(x)] + c [P_2(x) - P_1(x)]$
 - a, b, c are independent

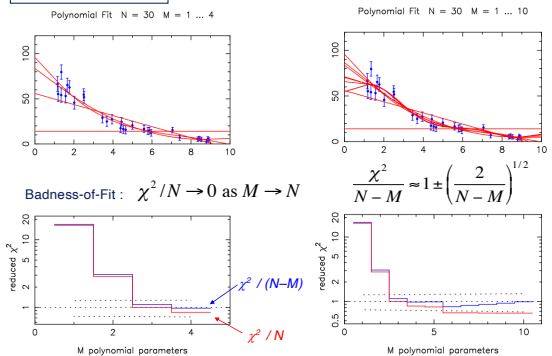


Note: every dataset has its own $1/\sigma^2$ weights, defines its own orthogonal polynomials.

219

Polynomial Fits

Mock data, true model is an exponential.



Badness-of-Fit: $\chi^2/N \rightarrow 0$ as $M \rightarrow N$

$$\frac{\chi^2}{N-M} \approx 1 \pm \left(\frac{2}{N-M} \right)^{1/2}$$

220

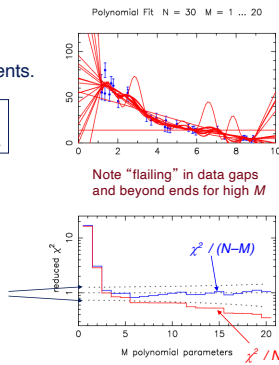
How many parameters to use ?

Fit $N = 30$ data points with $M = 1, 2, \dots, 20$ poly coefficients.

Higher $M =$ more flexible model.
Lower χ^2 , but less satisfactory fit.

χ^2_{\min} rejects $M = 1, 2$.
accepts $M = 3, 4, \dots$

$$\frac{\chi^2}{N-M} \approx 1 \pm \sqrt{\frac{2}{N-M}}$$



221

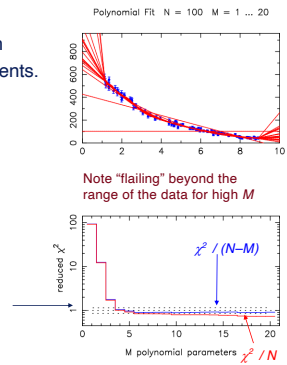
How many parameters to use ?

Fit $N = 100$ data points with $M = 1, 2, \dots, 20$ poly coefficients.

More data points and smaller error bars than before.

χ^2_{\min} rejects $M = 1, 2, 3$.
accepts $M = 4, 5, \dots$

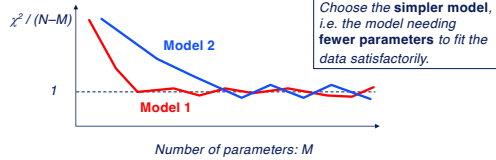
$$\frac{\chi^2}{N-M} \approx 1 \pm \sqrt{\frac{2}{N-M}}$$



222

Occam's Razor - "Keep it Simple"

- William of Occam (ca. 1286–1347):
"It is futile to do with more, what can be done with fewer"
or: "Keep it simple!"
- Fit 2 different models, 1 and 2, to the same data.
- Each model has $M=1, 2, \dots$ parameters, e.g. increasing numbers of polynomial coefficients.
- Prefer the simpler model.



223

Information Criteria: AIC, AICc, BIC

- Each parameter improves the fit: $-2 \ln(L)$ decreases.
- Include a penalty for each new parameter.
- Does the reduction in $-2 \ln(L)$ offset the penalty?

Minimise the AIC (or AICc or BIC) to choose the simplest model (needing the fewest parameters) that fits the data.

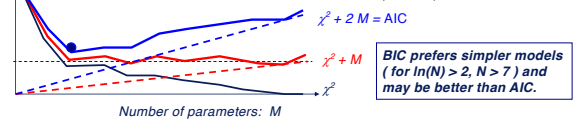
Data points: N Parameters: M Likelihood: L
Akaike Information Criterion: $AIC = -2 \ln(L) + 2M$

$$\text{Corrected AIC: } AICc = -2 \ln(L) + 2M \left(1 + \frac{M-1}{N} \right)$$

Bayesian Information Criterion: $BIC = -2 \ln(L) + \ln(N)M$

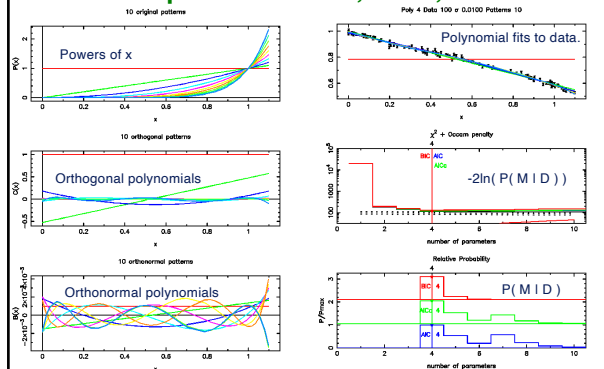
Likelihood: $P(\text{data} | M) \propto \exp(-\chi^2/2)$ Prior: $P(M) \propto \exp(-M)$

Posterior: $P(M | \text{data}) \propto P(\text{data} | M) P(M) \propto \exp(-AIC/2)$



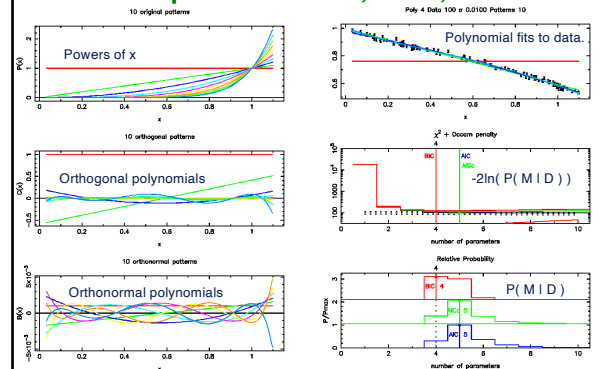
224

Comparison of AIC, AICc, BIC



225

Comparison of AIC, AICc, BIC



226

Fini -- ADA 11

227