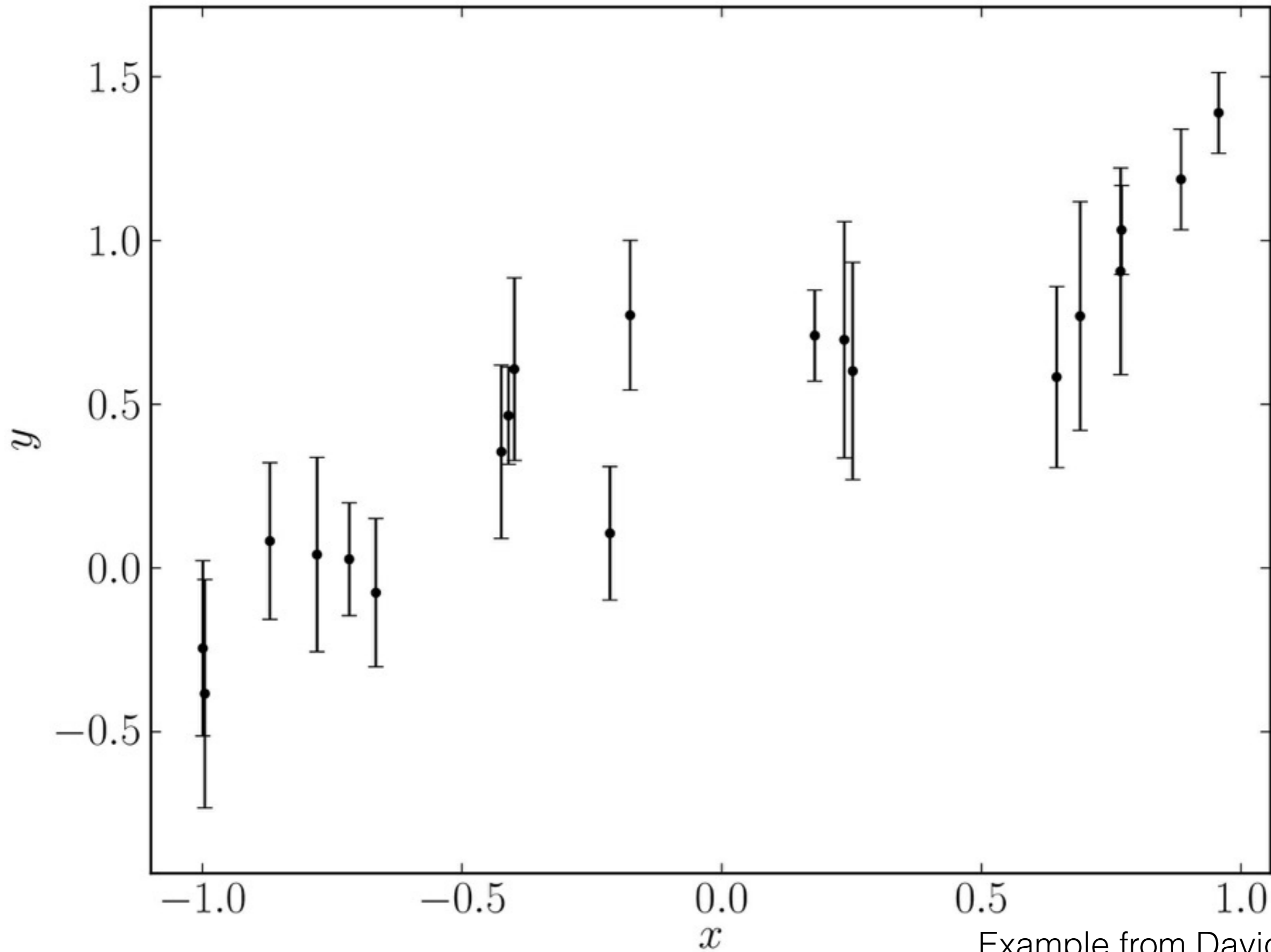


Statistics and Inference in Astrophysics

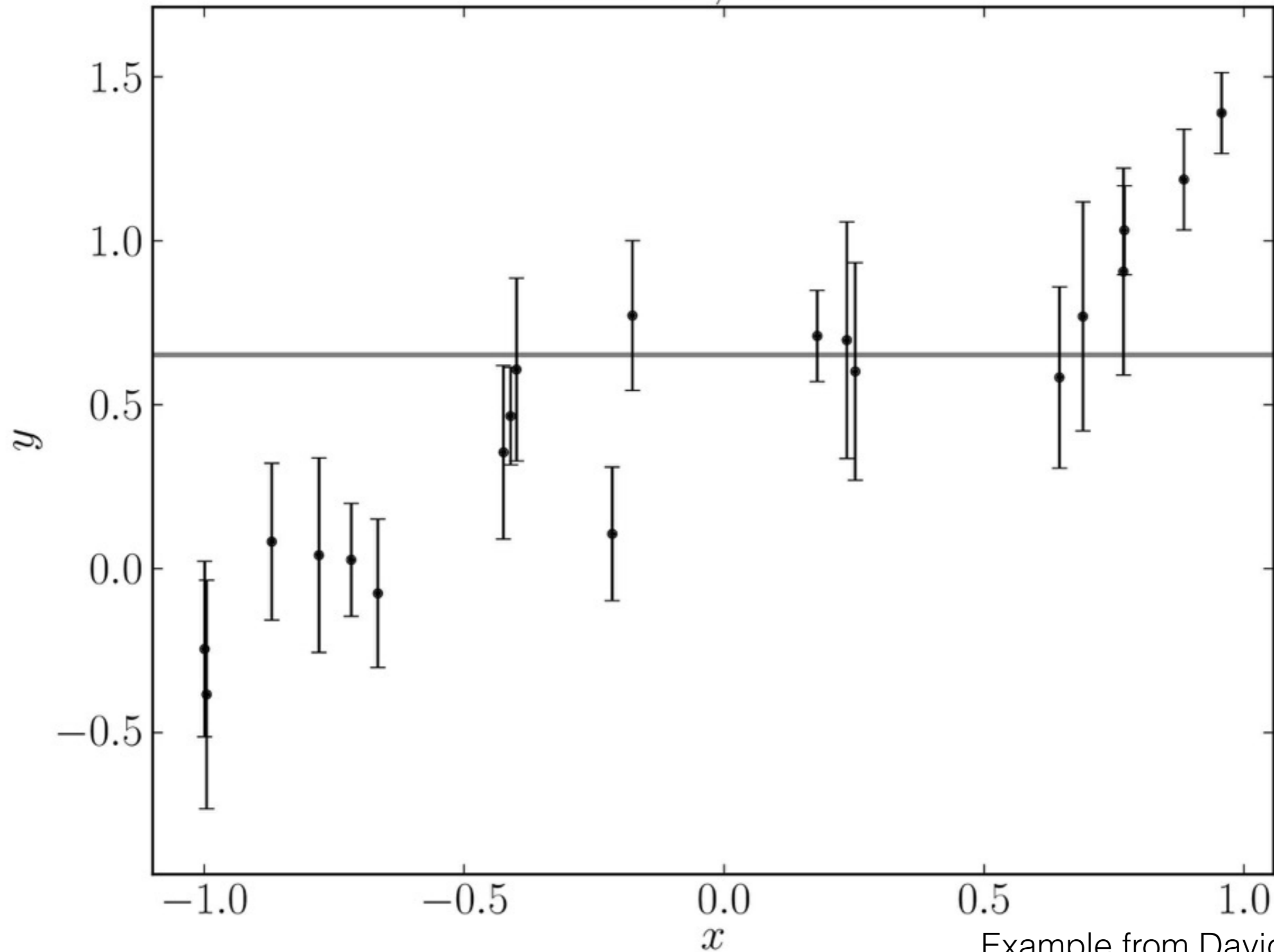
Today

- Goodness-of-fit, model selection, cross-validation
- Outliers, robust statistics



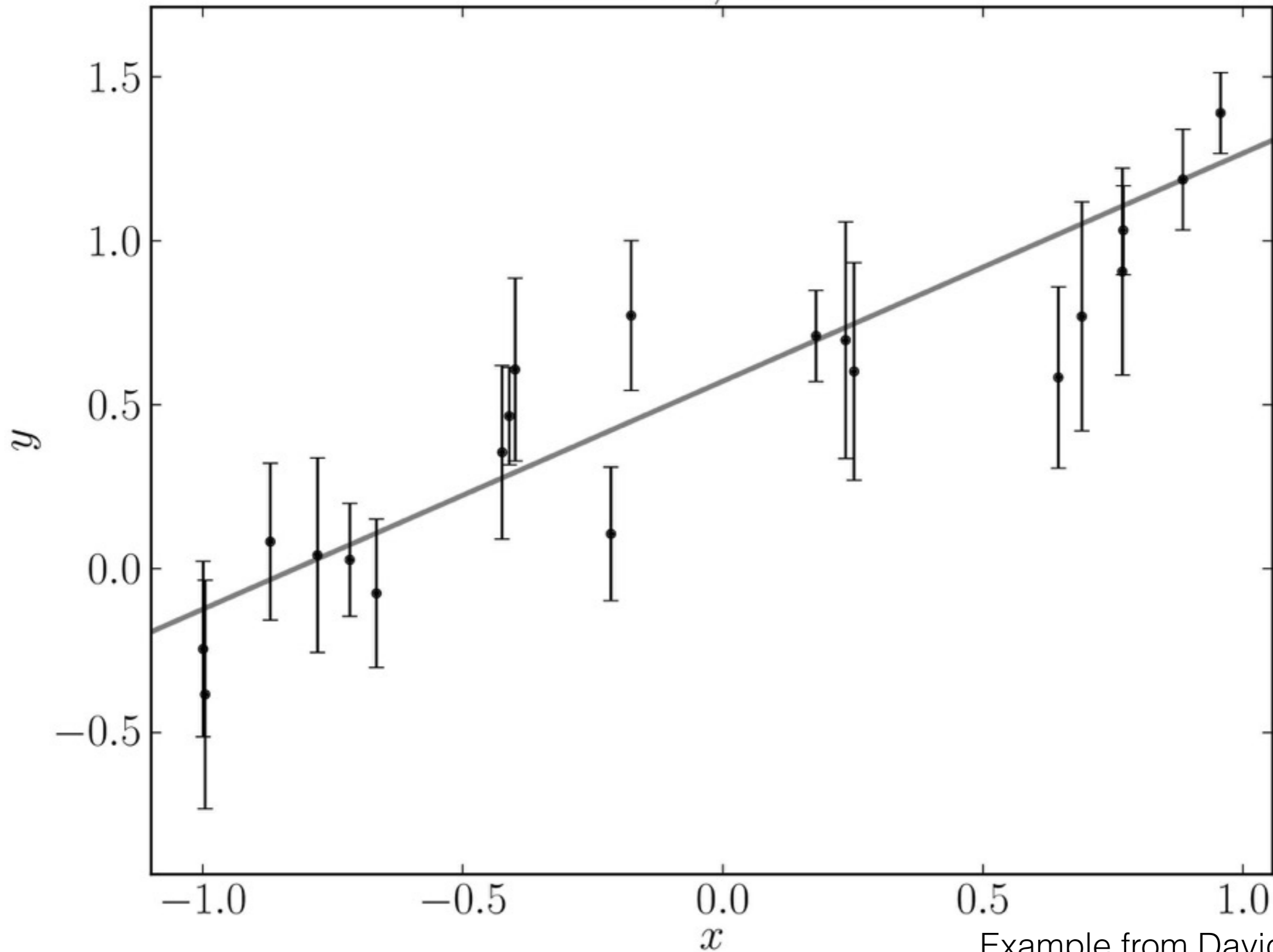
Example from David Hogg

order 0 ; $K = 1$



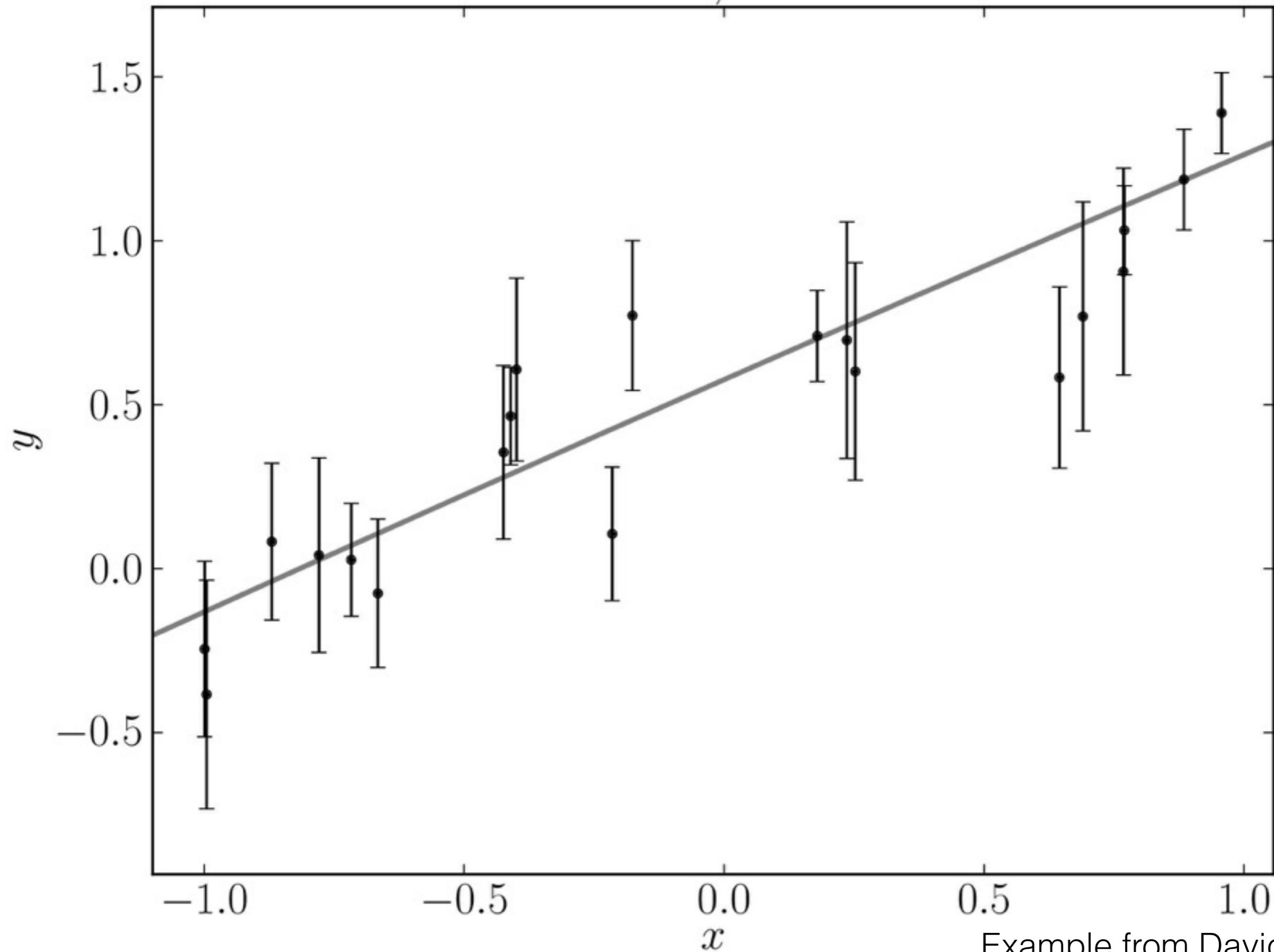
Example from David Hogg

order 1 ; $K = 2$



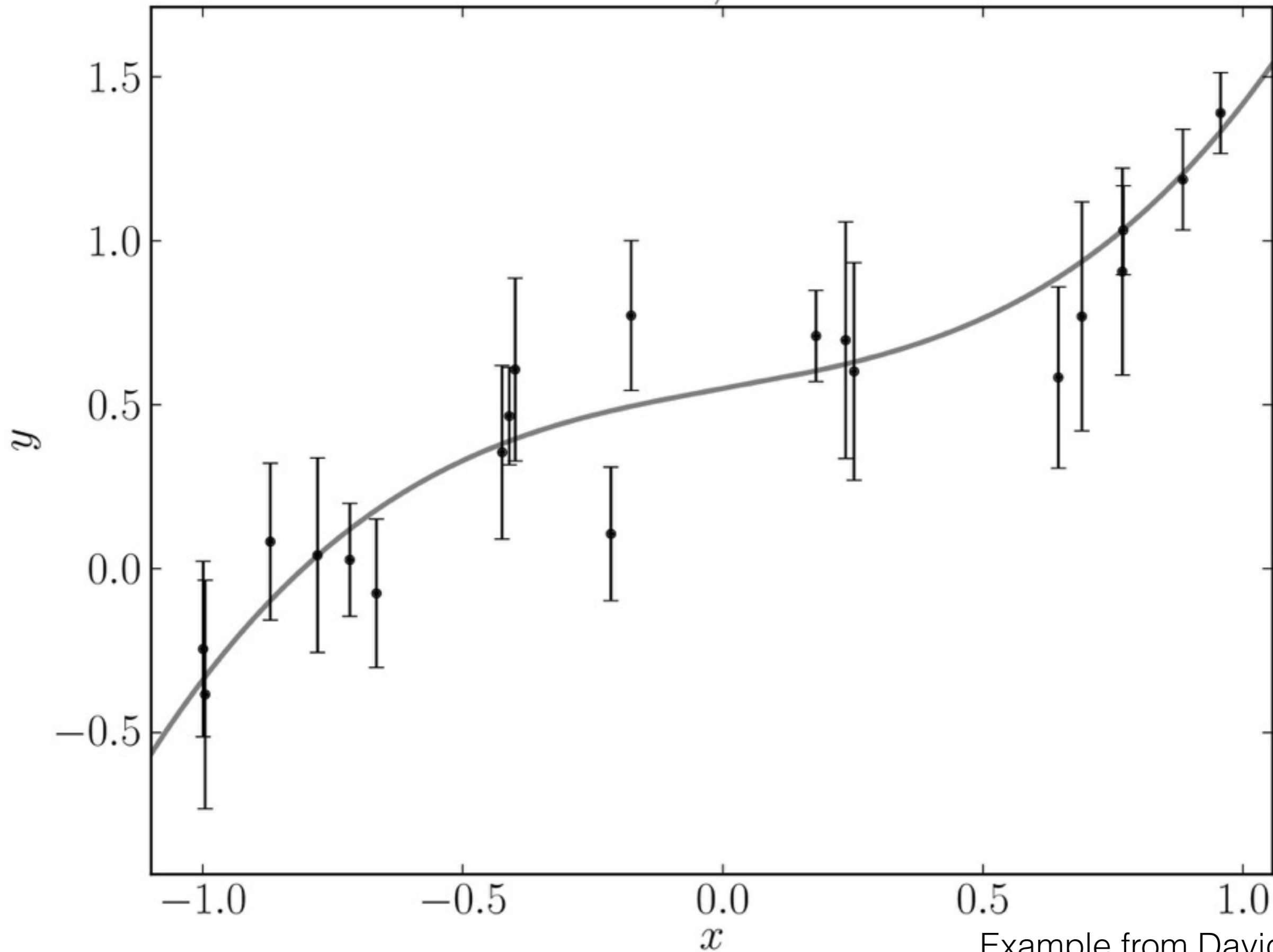
Example from David Hogg

order 2 ; $K = 3$



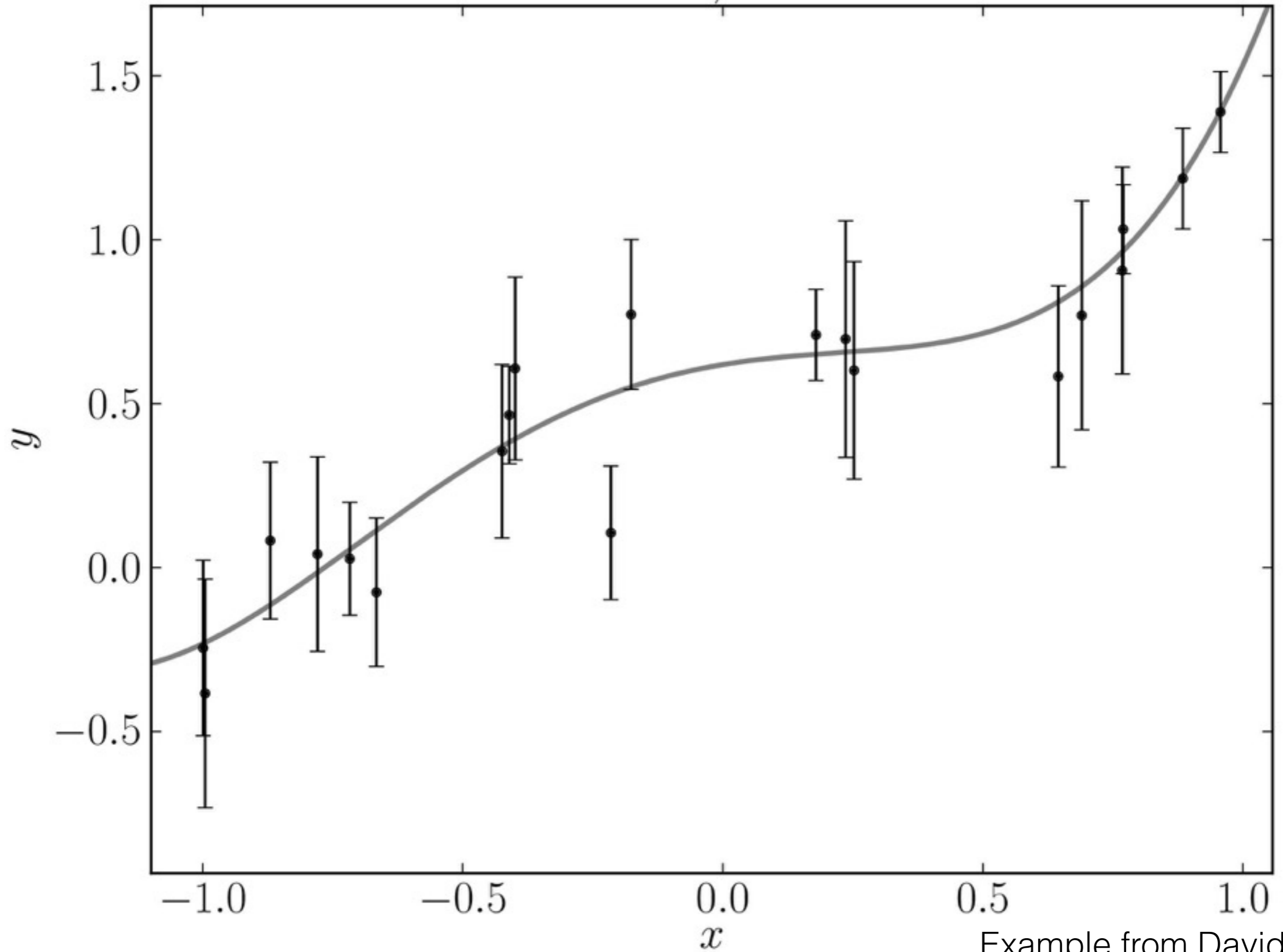
Example from David Hogg

order 3 ; $K = 4$



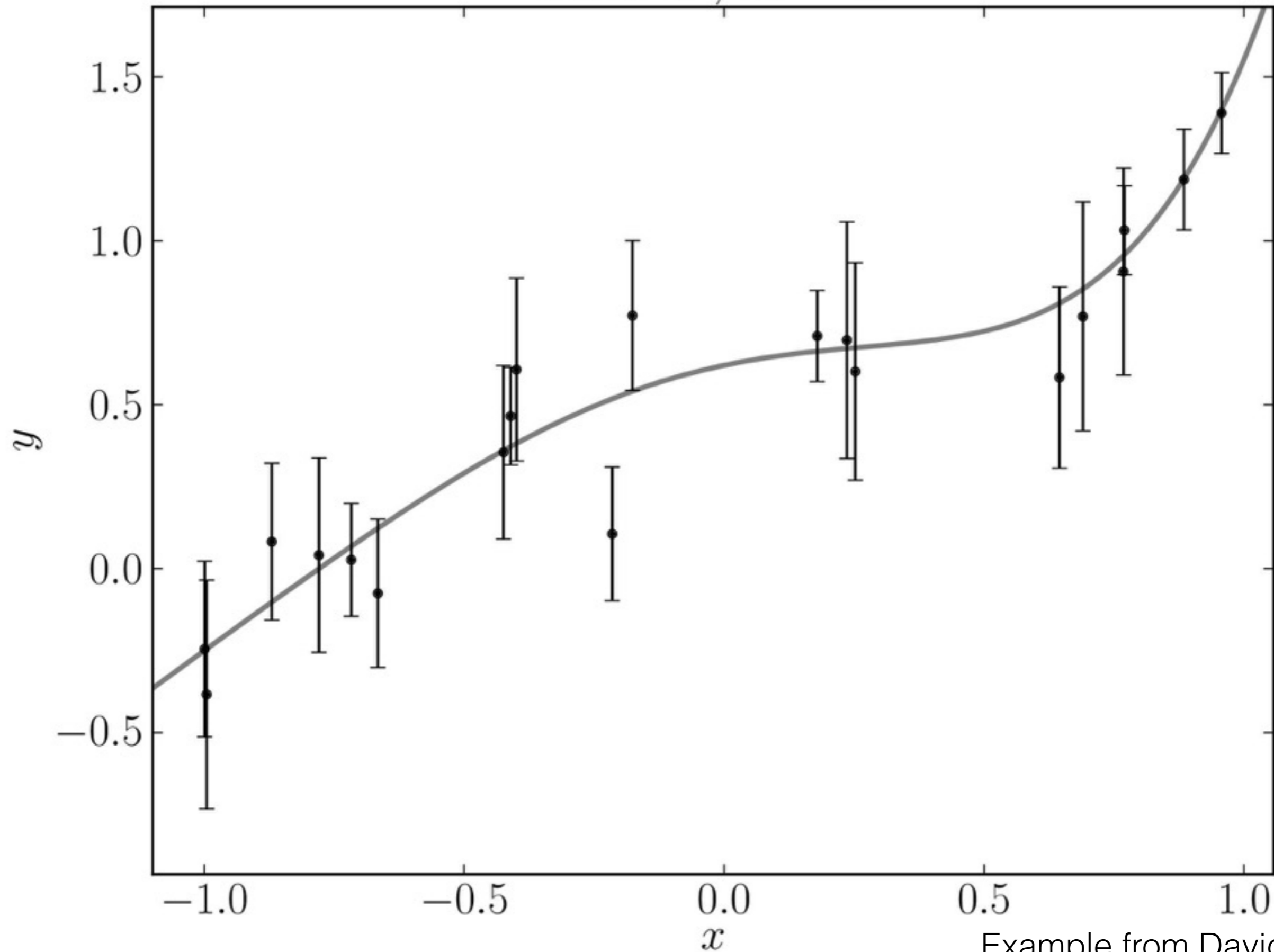
Example from David Hogg

order 4 ; $K = 5$



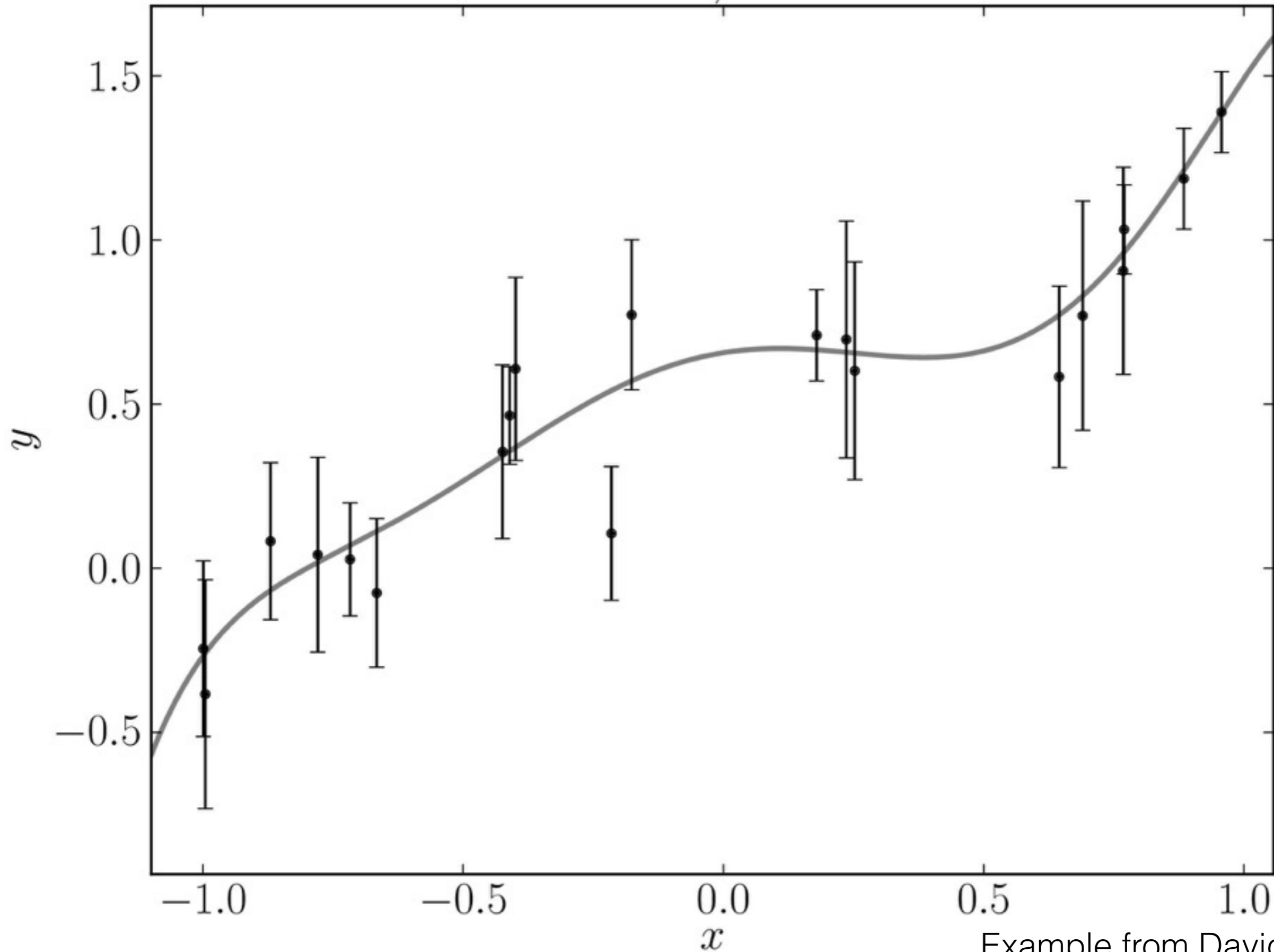
Example from David Hogg

order 5 ; $K = 6$



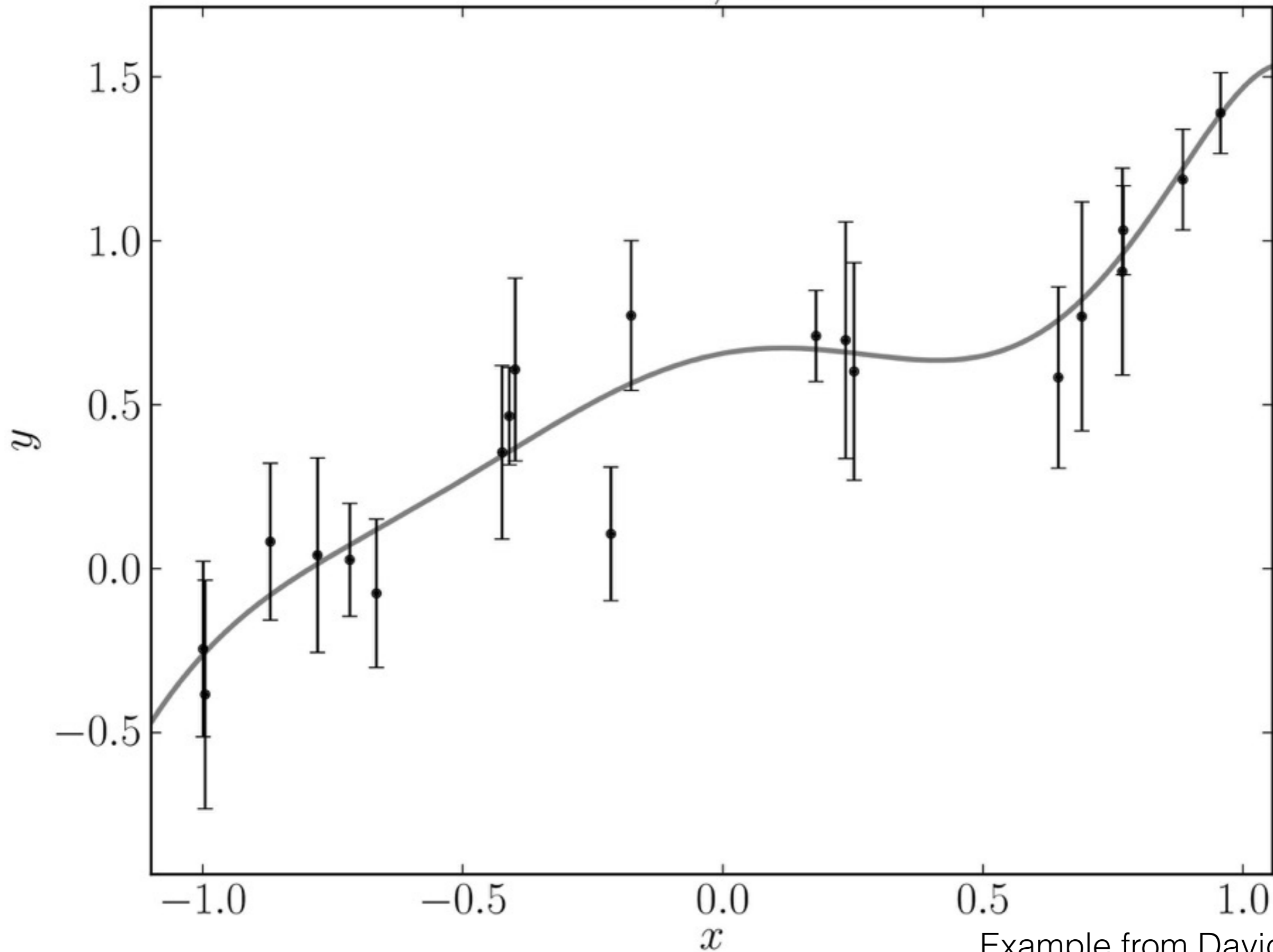
Example from David Hogg

order 6 ; $K = 7$



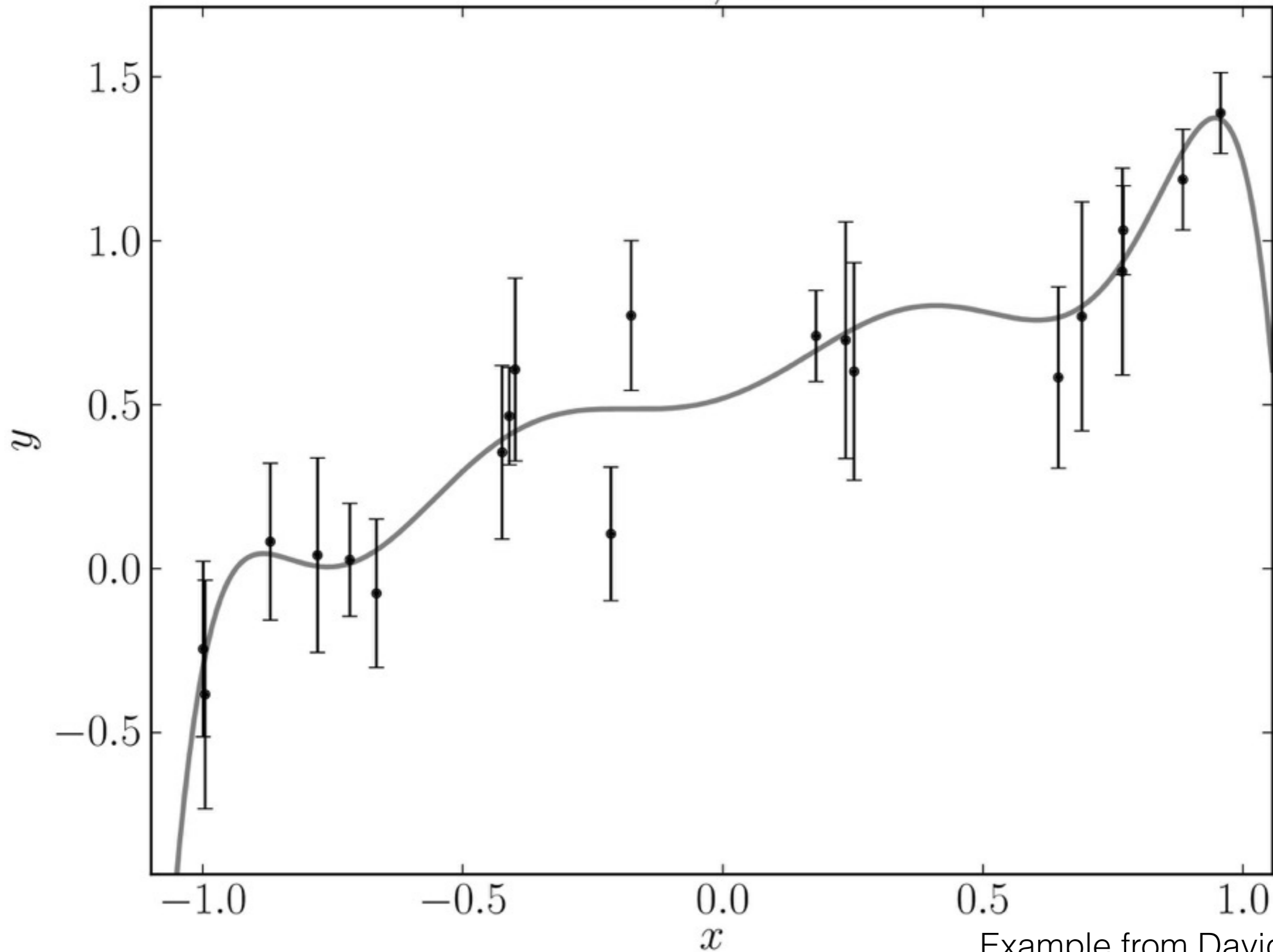
Example from David Hogg

order 7 ; $K = 8$



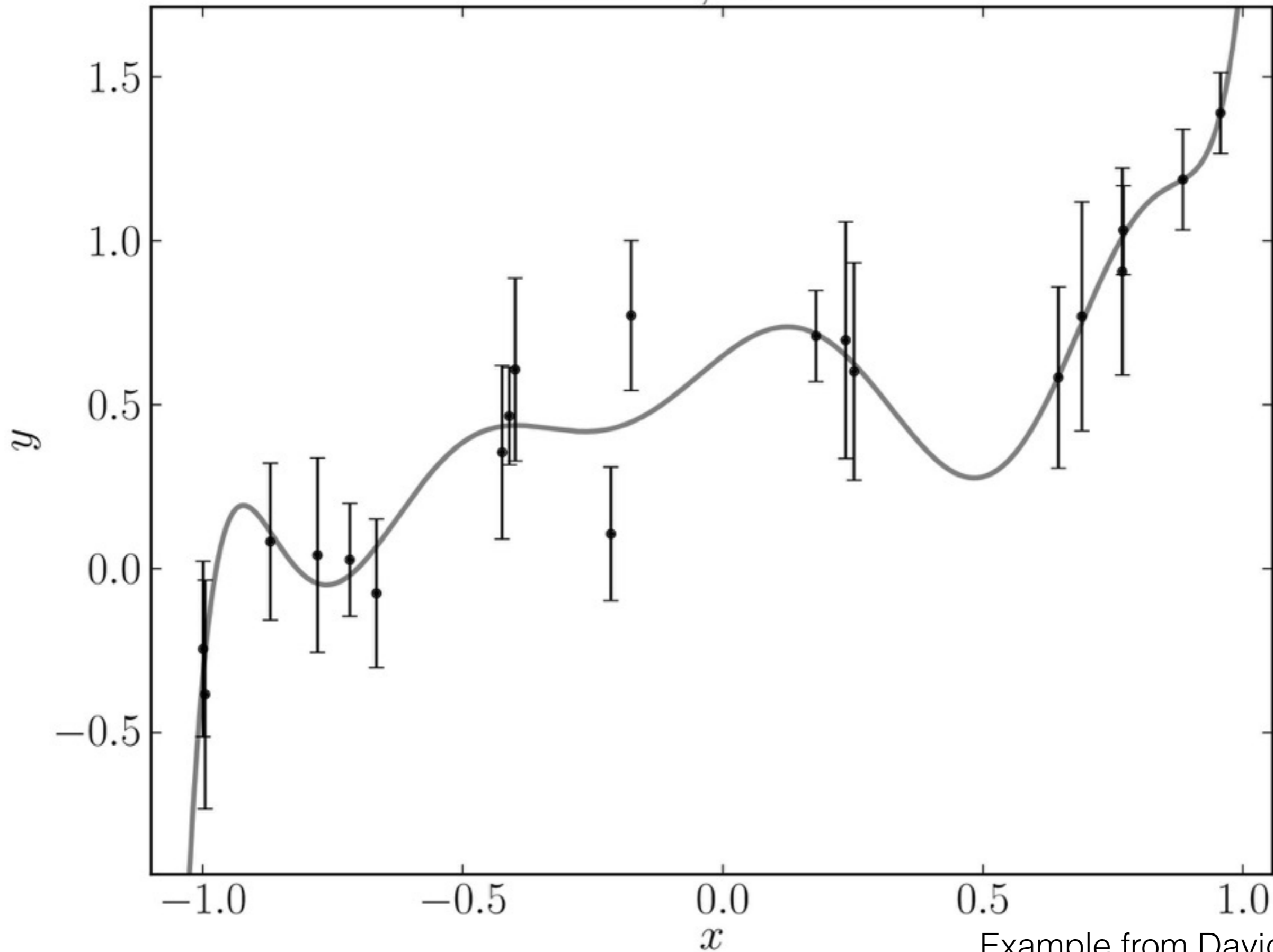
Example from David Hogg

order 8 ; $K = 9$



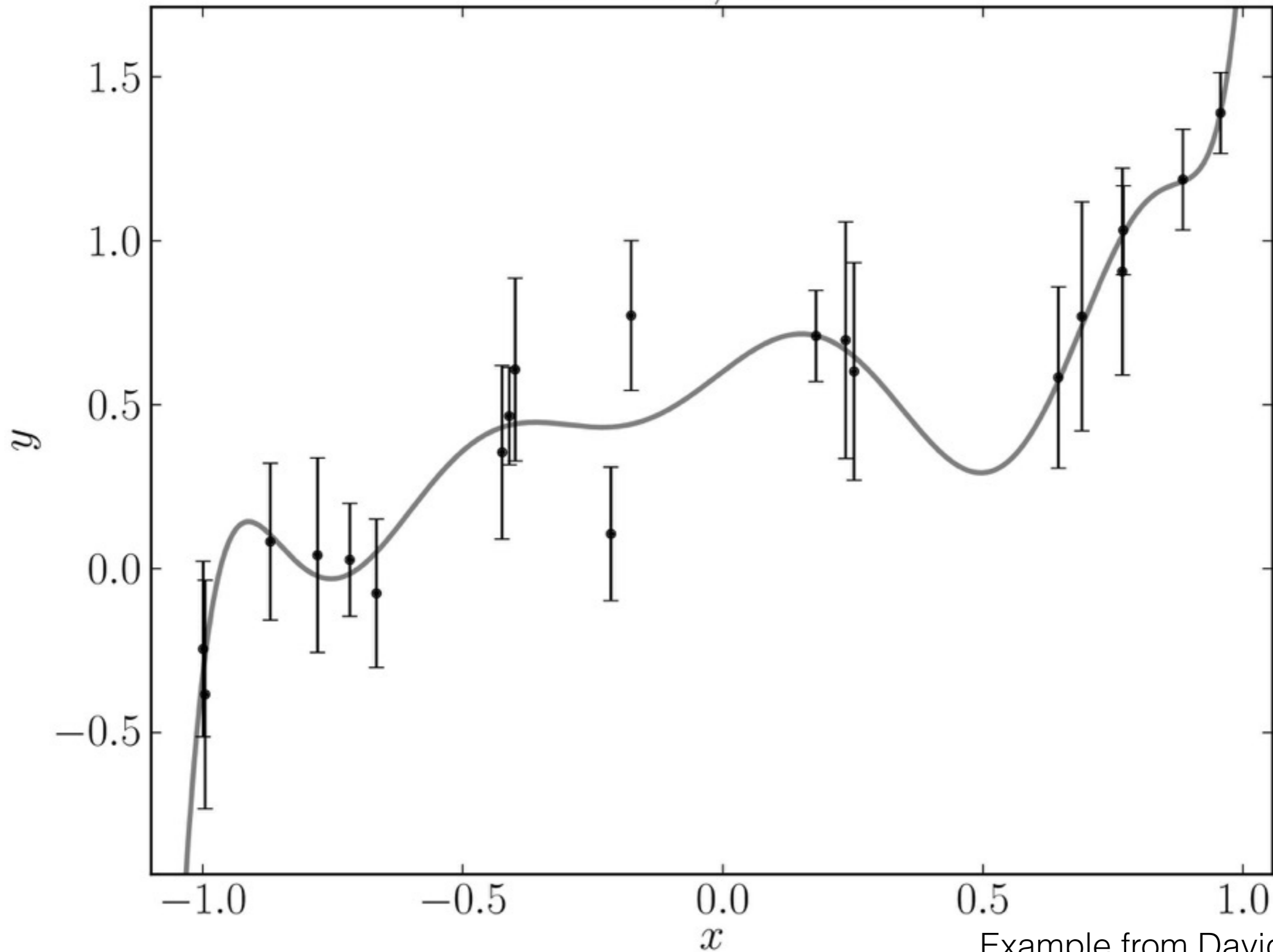
Example from David Hogg

order 9 ; $K = 10$



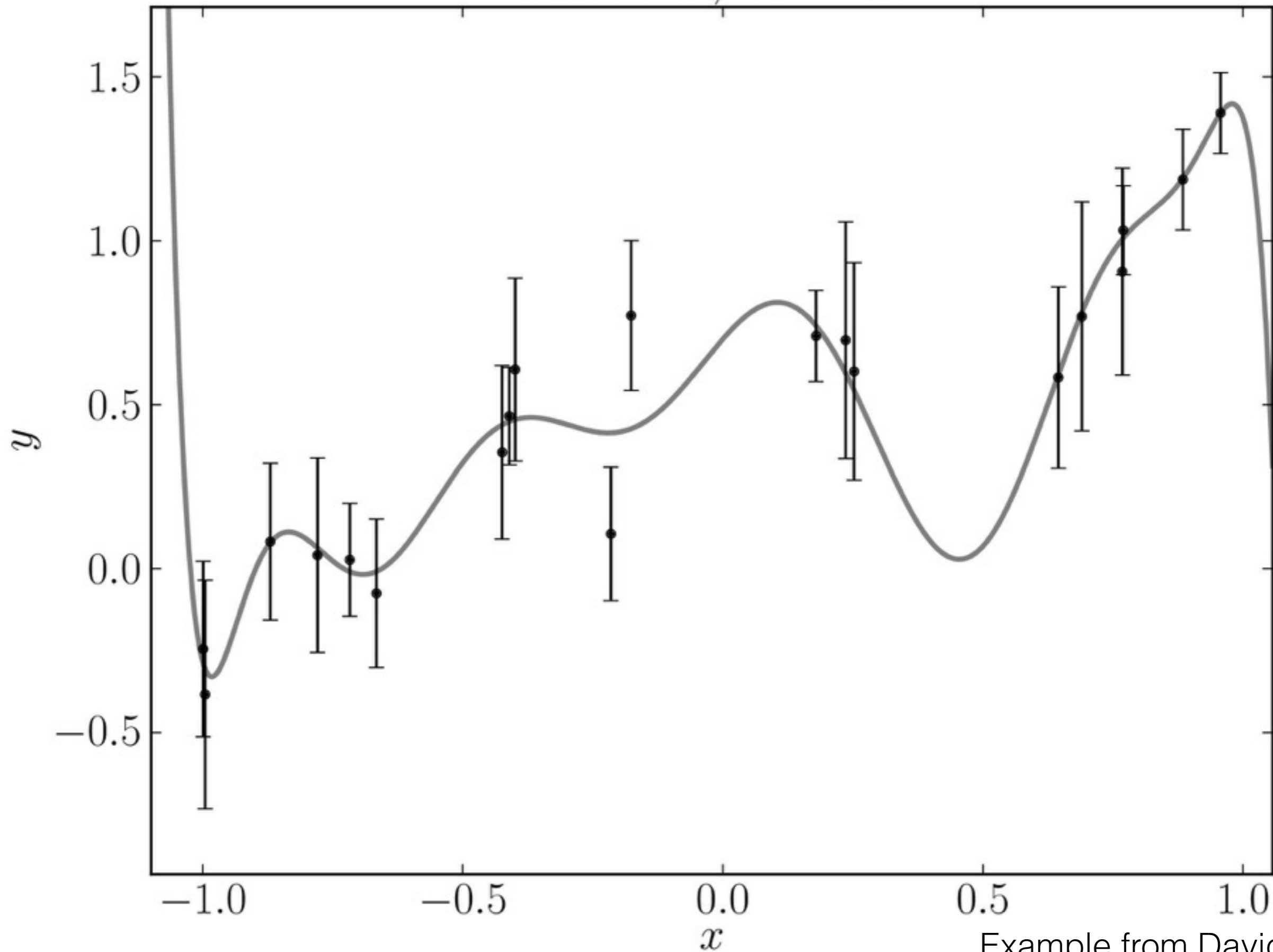
Example from David Hogg

order 10 ; $K = 11$



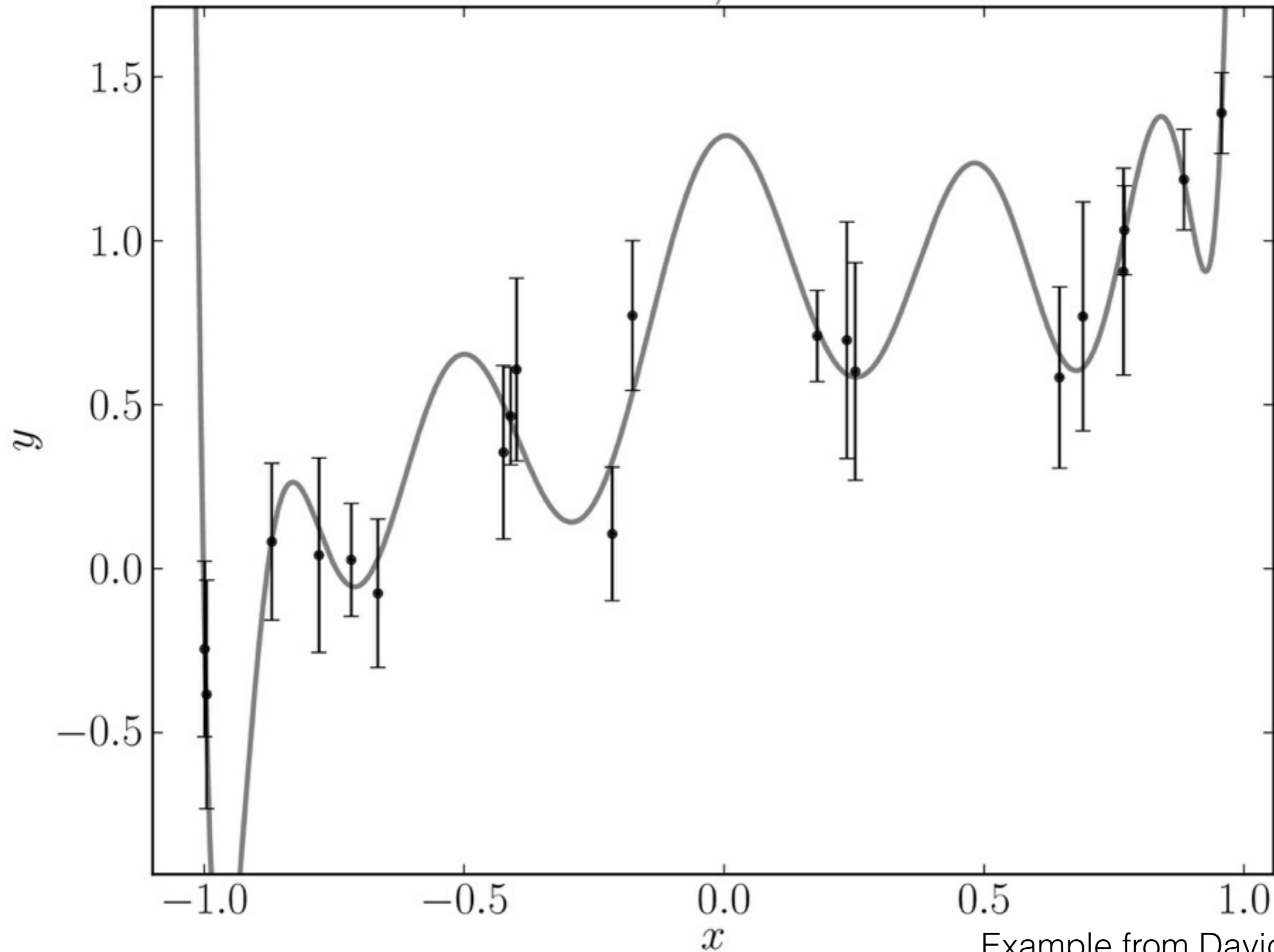
Example from David Hogg

order 11 ; $K = 12$



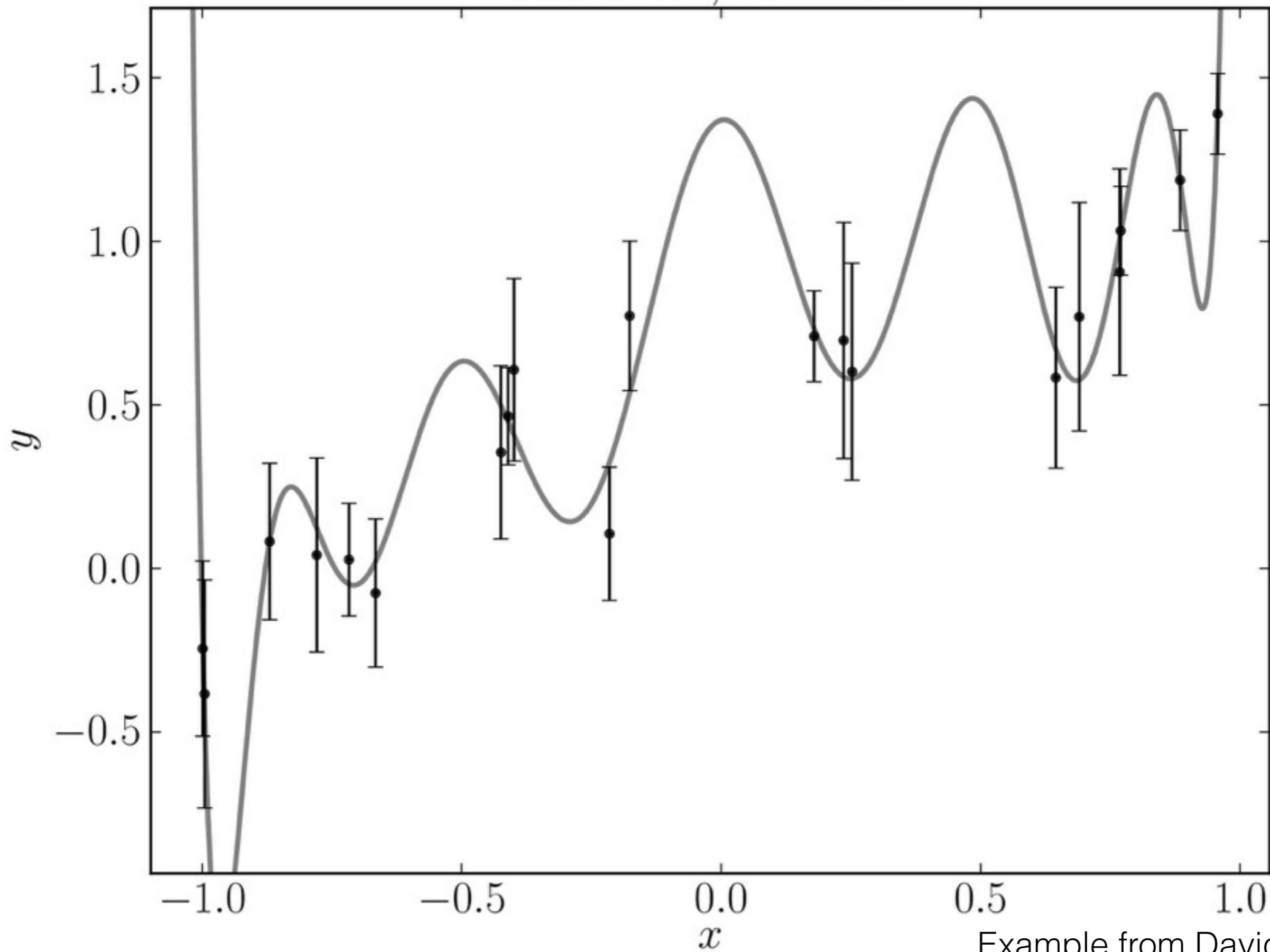
Example from David Hogg

order 12 ; $K = 13$



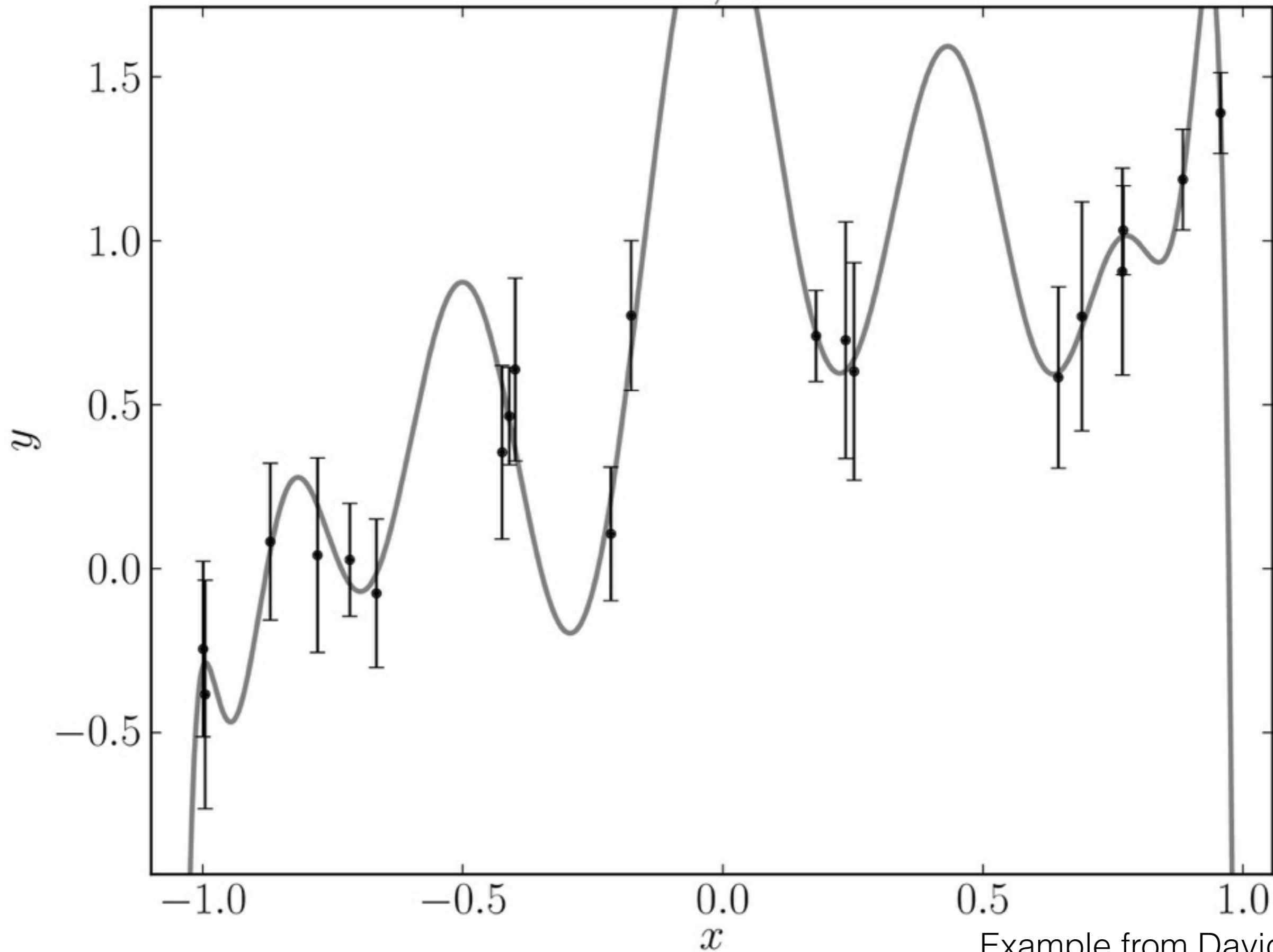
Example from David Hogg

order 13 ; $K = 14$



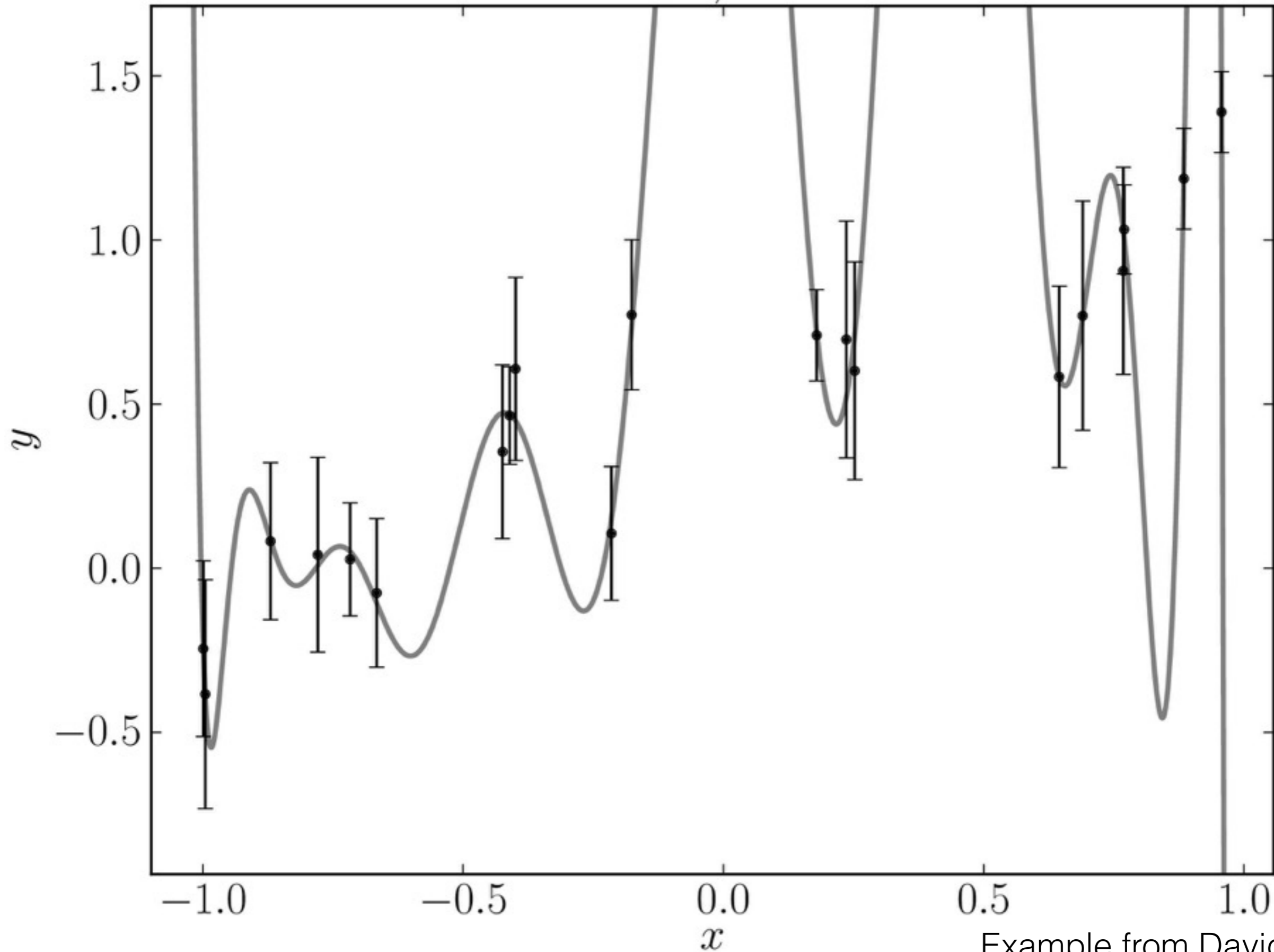
Example from David Hogg

order 14 ; $K = 15$

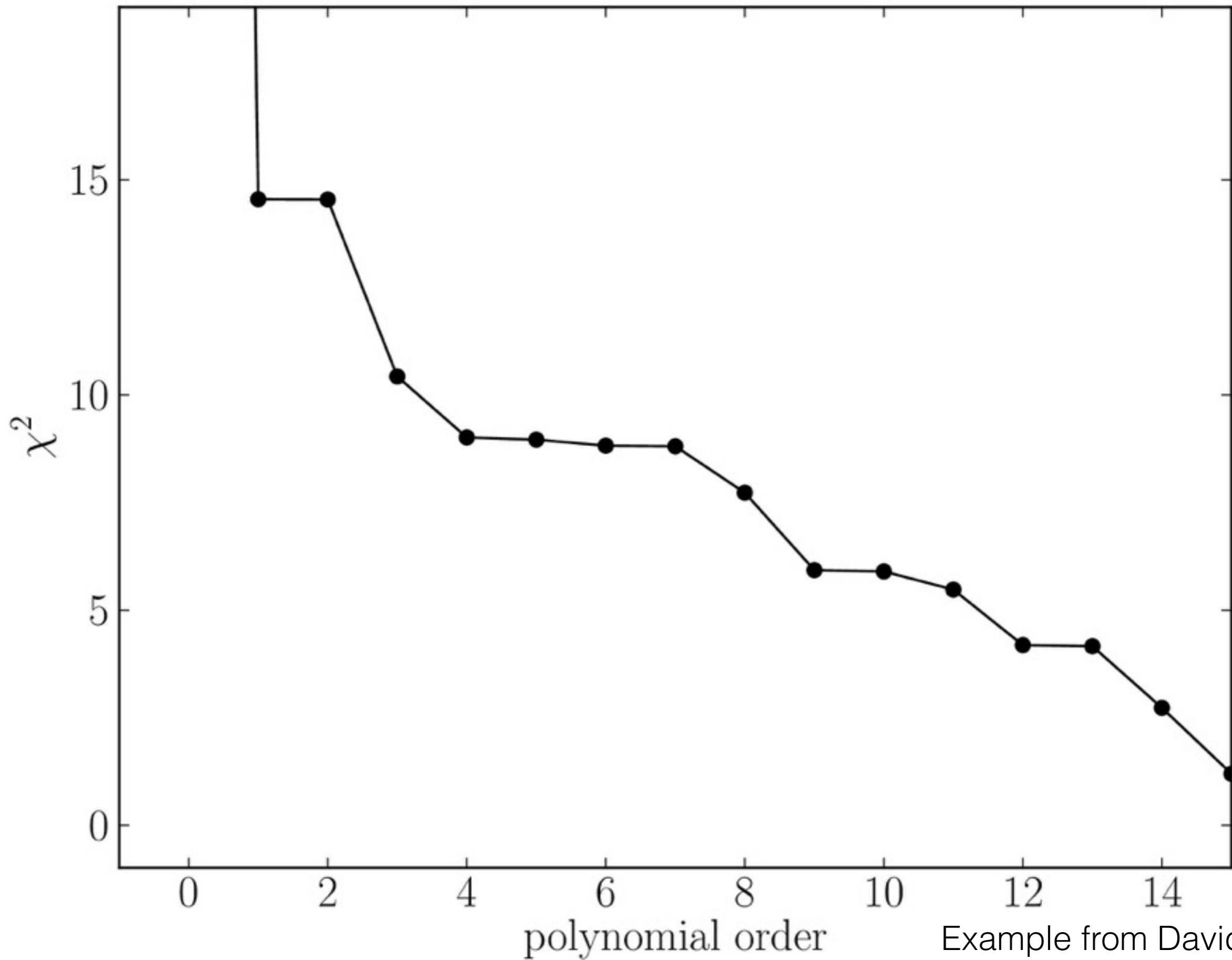


Example from David Hogg

order 15 ; $K = 16$



Example from David Hogg



Example from David Hogg

How do we decide which model is the best model?

- Chi-squared keeps improving as we increase the model's complexity
- But clear that we are overfitting!
- To determine the best model, need to figure out which model makes the *data the most probably* —> model selection
- Goodness-of-fit: is the data likely given the model?

Goodness-of-fit: General approach

- Given the model, simulate what the data would be like —> simulated data
- Compare this to the actual data that you have
- If the simulated data is very different from the actual data —> not a good fit!
- From simulated data, can reject that the data was generated by the model at X% confidence —> frequentist method at heart

Comparing simulated and actual data

- Large number of data summaries to choose from!
- Can look at plots, but for automated analysis require some low-dimensional summary
- Most popular: $\sum_i \chi_i^2$, $\sum_i |\chi_i|$, or look at distribution of χ_i

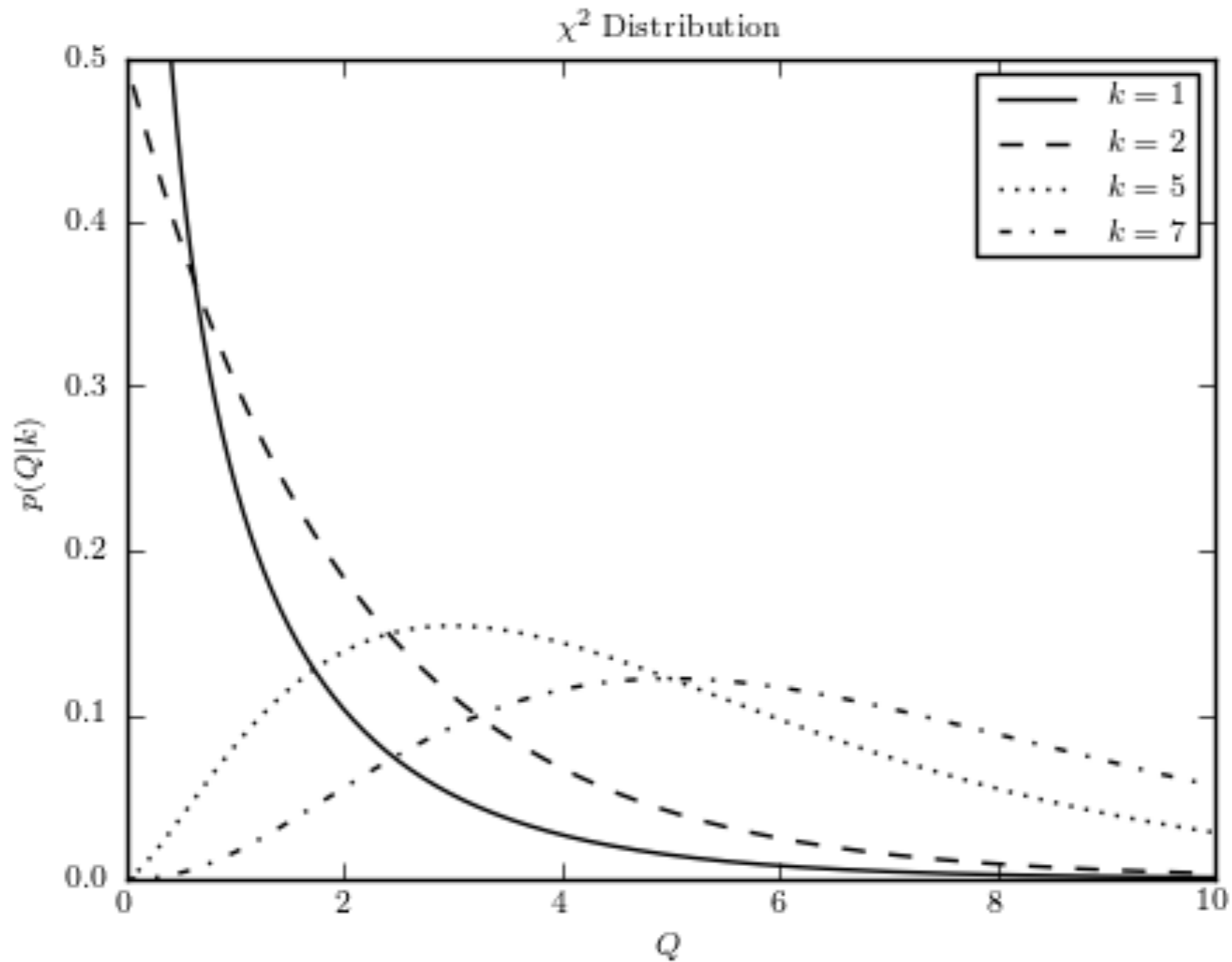
Chi-squared/ degree-of-freedom

- Most popular approach to goodness-of-fit uses chi-squared divided by the number of degrees of freedom
- Chi-squared here = $\sum_i \chi_i^2$
- Number of degrees of freedom = # data points - # of fit parameters
- Where does this come from? When does it apply?

Chi-squared distribution

- Distribution of sum of squares of k independent standard normal variables (those from $N(x|0,1)$)
- Form:
$$p(x|k) = \frac{1}{2^{k/2} \Gamma\left(\frac{k}{2}\right)} x^{k/2-1} e^{-x/2}$$
- Mean: k
- Variance: $2k$
- Basis for *chi-squared-per-degree-of-freedom* goodness-of-fit
- Central limit theorem: for $k \rightarrow \infty$, $p(x|k) \rightarrow N(x|k,2k)$

Chi-squared distribution



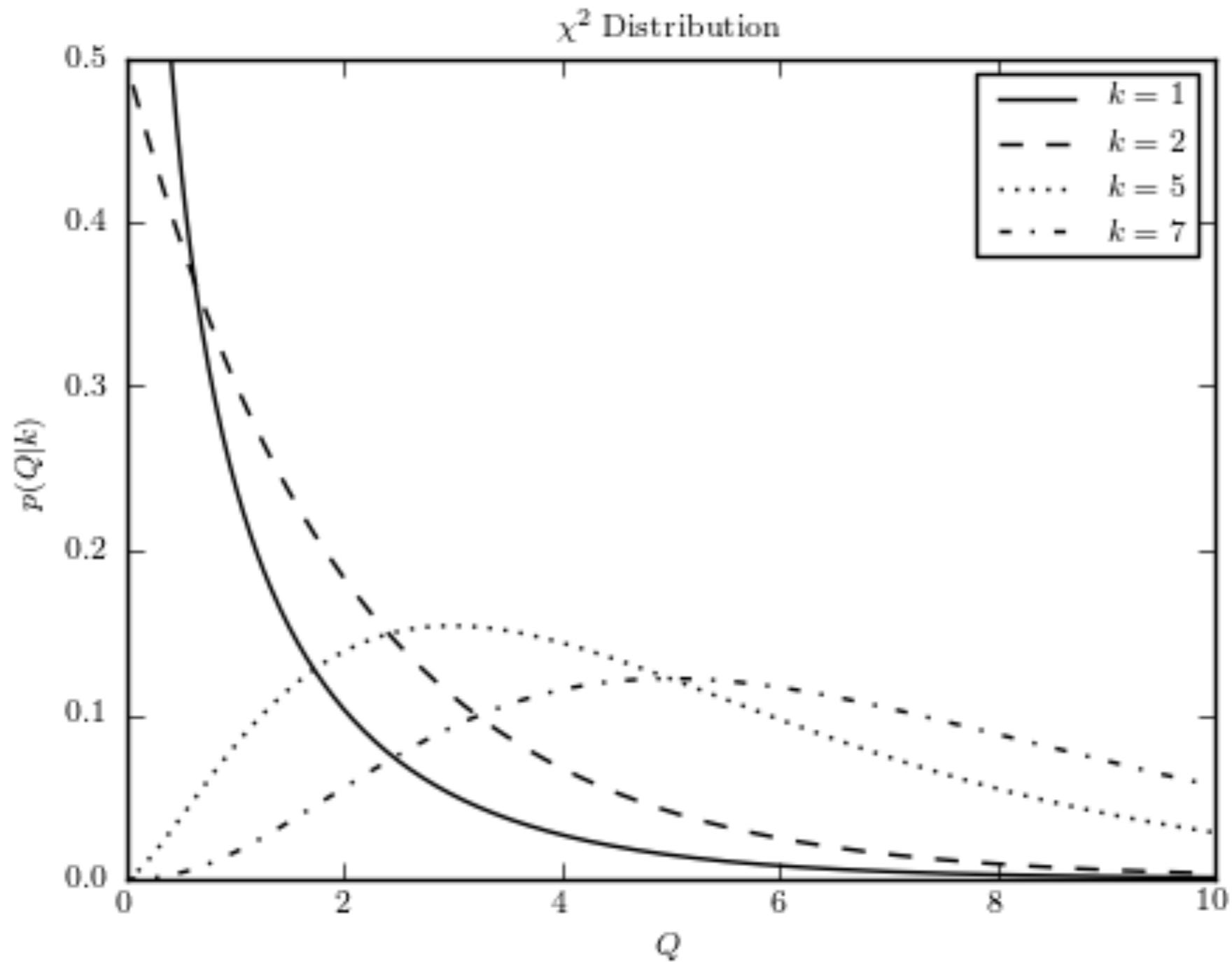
Chi-squared/ degree-of-freedom

- *If the likelihood is Gaussian: e.g., $p(y_i|m,b,x_i,\sigma_{y,i}) = N(y_i|mx_i+b,\sigma_{y,i}^2)$*
- *Then we have: $-2 \ln L = \sum_i \chi_i^2$, e.g., $= \sum_i [(y_i-mx_i-b)/\sigma_{y,i}]^2$*
- *and $\chi_i \sim N(0, 1)$*
- *Therefore, $\sum_i \chi_i^2$ is distributed as a chi-squared distribution*
- *If we have fit K parameters to N data, only $N-K$ of these χ_i^2 are independent $\rightarrow \sum_i \chi_i^2$ is chi-squared distributed with $N-K$ degrees of freedom*

Chi-squared distribution with $N-K$ degrees of freedom

- Mean = $N-K$ = dof
- Variance = $2(N-K)$ = 2 dof
- Central limit theorem:
for $N-K \rightarrow \infty$, $p(x|N-K) \rightarrow N(x| N-K, 2[N-K])$
- Therefore, expected value of $\chi^2/\text{dof} \sim 1$
- But really should be comparing χ^2 to dof with typical scatter $\sqrt{2\text{dof}}$

Chi-squared distribution



Chi-squared/ degree-of-freedom

- Assumptions one more time!
- Likelihood is Gaussian \rightarrow for Gaussian uncertainties means that the model must be *linear* (e.g., polynomial)
- Must *believe* the uncertainties
- #dof must be large \rightarrow large data limit
- Almost never directly applies in practice! But for well-constrained parameters, any model space approx. linear near the best-fit \rightarrow widespread use of χ^2/dof

Chi-squared/ degree-of-freedom for non-linear / non-Gaussian models

- If the likelihood (data uncertainty) is not Gaussian or the model is not linear, χ^2/dof does not technically apply (except in the limit discussed on the previous slide)
- So could just directly simulate the data (e.g., linear-fit):
 1. For best fit model parameters (m,b)
 2. Simulate data: $y = mx + b$
 3. Draw random uncertainties and add them to y
 4. Compute χ^2
 5. Repeat to form $p(\chi^2)$
 6. Where does χ^2 for the actual data lie in this distribution?
- Similar with other summaries (don't need to use χ^2)

Model selection using chi-squared/ degree-of-freedom

- χ^2/dof can be used to select the best model among *linear* models
- Idea is that when overfitting, χ^2 will be suspiciously close to zero
- When underfitting, χ^2 will be large
- Best model makes the data most likely \rightarrow peak of χ^2 distribution $\rightarrow \chi^2/\text{dof} \sim 1$

Detour: difference between Delta χ^2 and χ^2/dof

- χ^2 comes up in model fitting and goodness-of-fit
- $-2 \ln L = \chi^2 = \sum_i \chi_i^2$
- Find best fit, can compute $\Delta\chi^2 = \chi^2 - \chi_{\min}^2$
- With uniform prior: $\ln p(m,b|\text{data}) = -\chi^2/2 \sim \Delta\chi^2 = \chi^2 - \chi_{\min}^2$
- Again, chi-squared distribution, but now with K degrees of freedom

Detour: difference between Delta χ^2 and χ^2/dof

- Suppose you have 1 parameter (e.g., fit constant $y = b$ rather than $y = m x + b$)
- $p(b|\text{data}) = \text{Chi}^2(\Delta\chi^2, 1 \text{ dof})$
- 68% confidence limit on b : that value for which $\Delta\chi^2 = 1$

```
In [7]: from scipy import stats
```

```
In [8]: stats.chi2.ppf(0.68, 1)
```

```
Out[8]: 0.98894648147802289
```

Detour: difference between Delta χ^2 and χ^2/dof

- Suppose you have 2 parameters (e.g., linear fit $y = m x + b$)
- $p(m,b|\text{data}) = \text{Chi}^2(\Delta\chi^2, 2 \text{ dof})$
- 68% confidence limit on (m,b): that ellipse for which $\Delta\chi^2 = 2.3$

```
In [11]: from scipy import stats
```

```
In [12]: stats.chi2.ppf(0.68,2)
```

```
Out[12]: 2.2788685663767296
```

Detour: difference between Delta χ^2 and χ^2/dof

$\Delta\chi^2$ as a Function of Confidence Level p and Number of Parameters of Interest ν						
p	ν					
	1	2	3	4	5	6
68.27%	1.00	2.30	3.53	4.72	5.89	7.04
90%	2.71	4.61	6.25	7.78	9.24	10.6
95.45%	4.00	6.18	8.02	9.72	11.3	12.8
99%	6.63	9.21	11.3	13.3	15.1	16.8
99.73%	9.00	11.8	14.2	16.3	18.2	20.1
99.99%	15.1	18.4	21.1	23.5	25.7	27.9

Numerical recipes 3rd edition

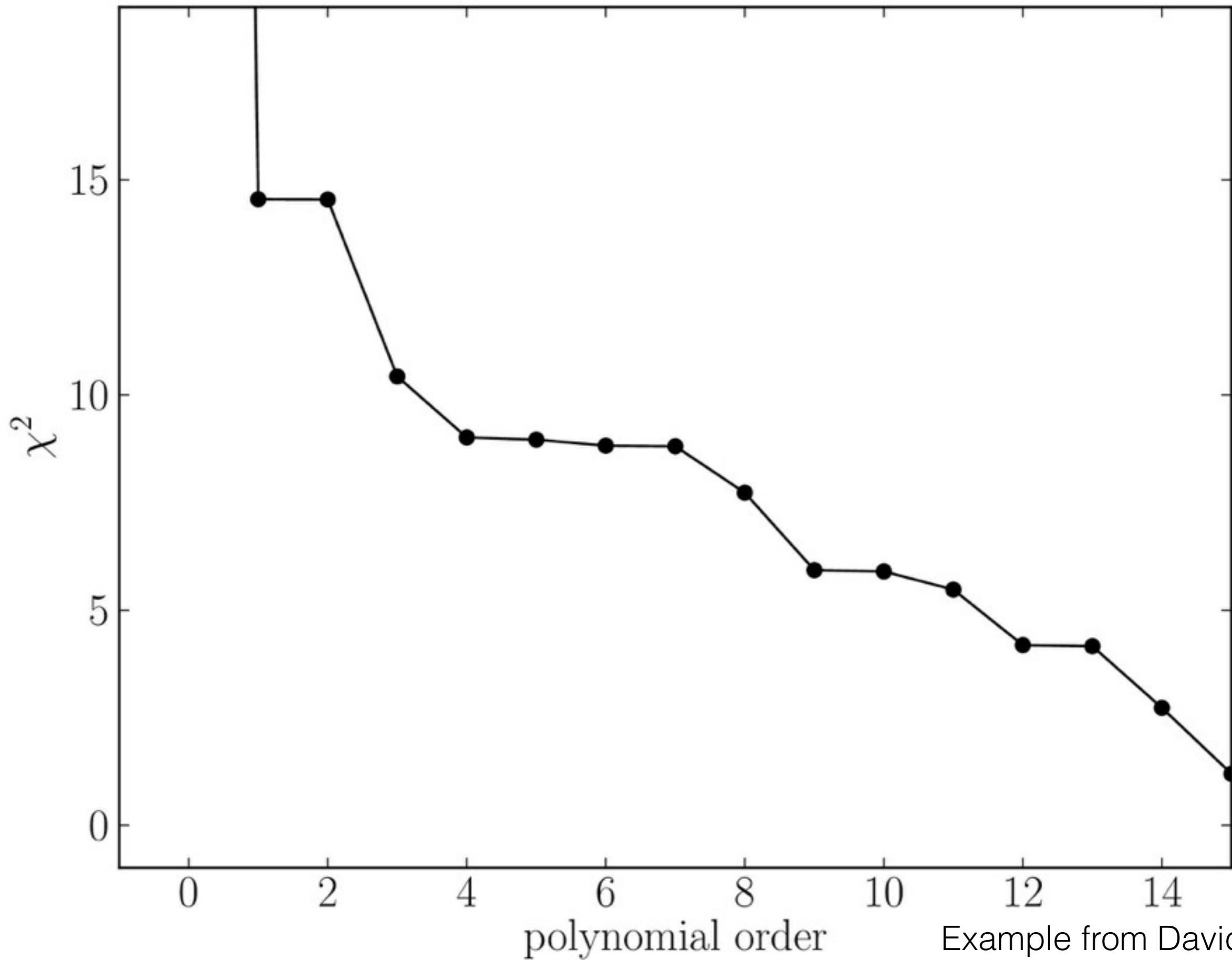
Detour: difference between Delta χ^2 and χ^2/dof

- Delta χ^2 used for finding uncertainty limits
- ***Don't*** use Delta χ^2 / dof for this! limits will always be much wider
- χ^2 / dof for model selection and goodness-of-fit

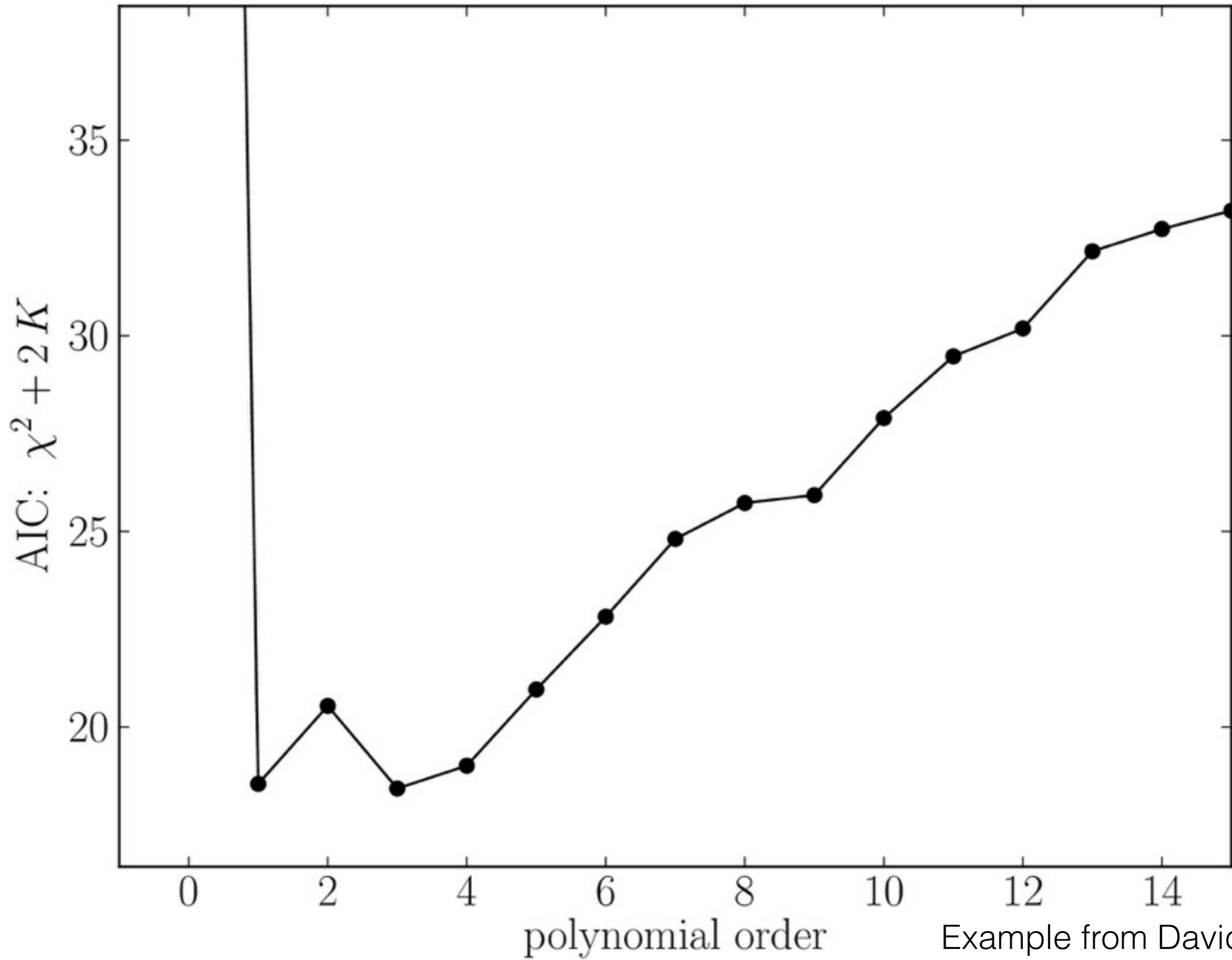
Model selection with AIC

- AIC = Akaike Information Criterion
- $AIC = -2 \ln L + 2 K = \chi^2 + 2 K$
- Delta AIC = Asymptotically the amount of information lost when using worse model G1 than better model G2
- Corrections for finite sample sizes depend on model;
for 1D, linear, Gaussian model

$$2K(K+1)/(N-K-1)$$



Example from David Hogg



Example from David Hogg

Bayesian model selection

- Bayesian methods are *very bad* at goodness-of-fit
- Application of Bayes's theorem only allows to distinguish between 2 different models, no real concept of 'good model'
- Bayes's theorem for models: model1 = linear, 2 = quadratic

$$p(\text{linear}|\text{data}) \sim p(\text{data}|\text{linear}) p(\text{linear})$$

$$p(\text{quadratic}|\text{data}) \sim p(\text{data}|\text{quadratic}) p(\text{quadratic})$$

- Likelihoods in these equations are marginalized over parameters of each model (c = quadratic coeff.) \rightarrow marginalized likelihood

$$p(\text{data}|\text{linear}) = \int dm \, db \, p(\text{data}|m,b) p(m,b|\text{linear})$$

$$p(\text{data}|\text{quadratic}) = \int dc \, dm \, db \, p(\text{data}|m,b,c) p(m,b,c|\text{quadratic})$$

These were in the denominator of Bayes's theorem for $p(m,b|\text{linear})$ before!

Bayesian model selection

- Can then compute *odds ratio*:

$$\text{odds ratio} = \frac{p(\text{linear}|\text{data})}{p(\text{quadratic}|\text{data})}$$

and select the model with the highest odds

- Requires prior over models [$p(\text{linear})$ and $p(\text{quadratic})$]

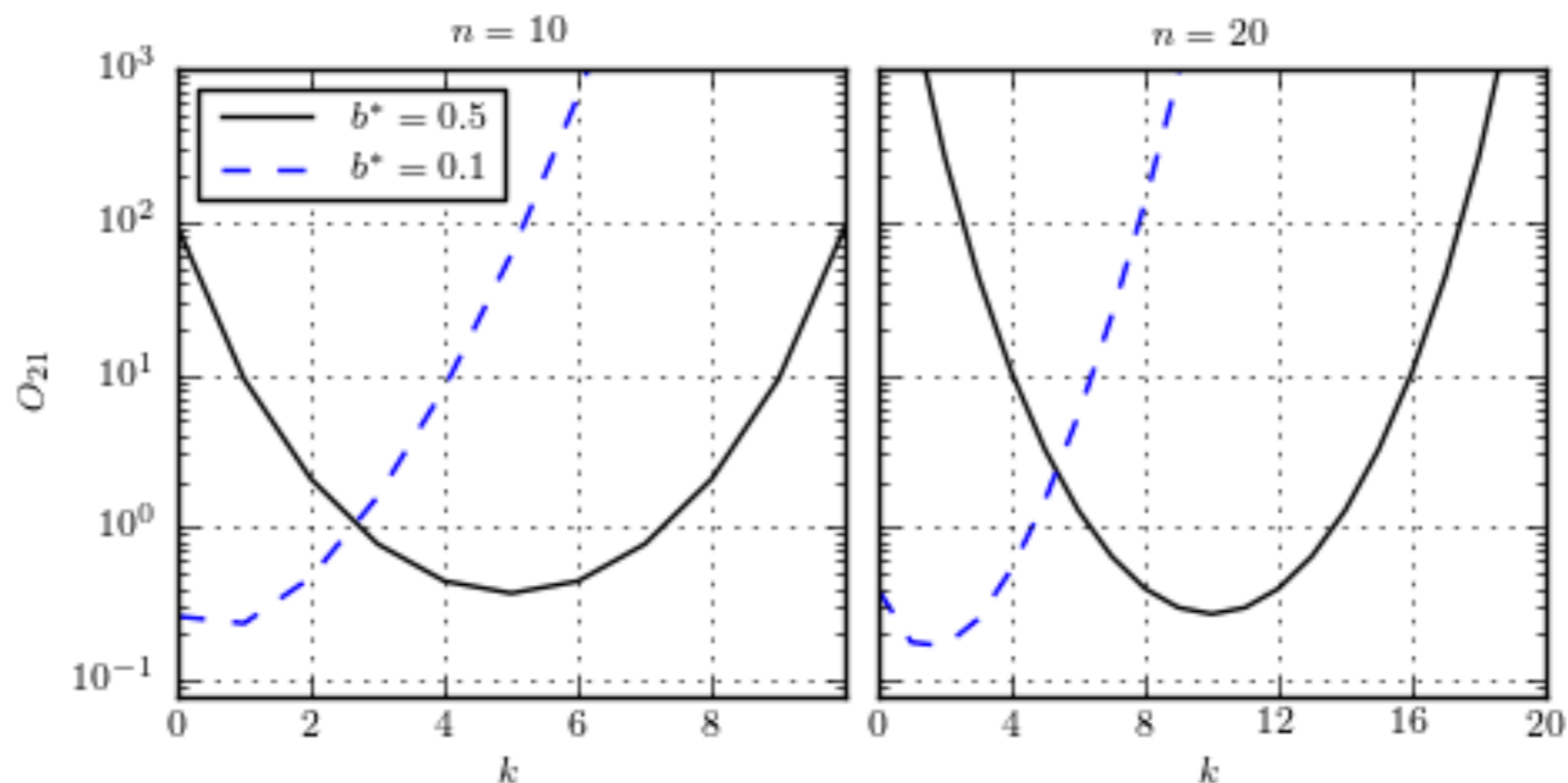
$$\text{odds ratio} = \frac{p(\text{data}|\text{linear})}{p(\text{data}|\text{quadratic})} \times \frac{p(\text{linear})}{p(\text{quadratic})}$$

$$= \text{Bayes-factor} \times \text{prior-ratio}$$

- If not strong preference, select based on Bayes-factor

Bayesian model selection: Example: is a coin fair?

- Flip coin n times, get k heads \rightarrow is it fair?
- Bayesian needs alternative model! Make that model: constant probability for heads that is unknown
- So two models: $p=0.5$ and $p=\text{unknown}$



Bayesian model selection: How to compute the evidence

- Evidence = $p(\text{data}|\text{model1}) = \int d\text{params } p(\text{data}|\text{params}) p(\text{params}|\text{model1})$
- Nested sampling: MCMC technique that returns the evidence
- If the posterior is close to Gaussian: Laplace approximation around $\max P^*(x)$ with covariance matrix **A**:

$$P^*(\mathbf{x}_0) \sqrt{\frac{(2\pi)^K}{\det \mathbf{A}}}$$

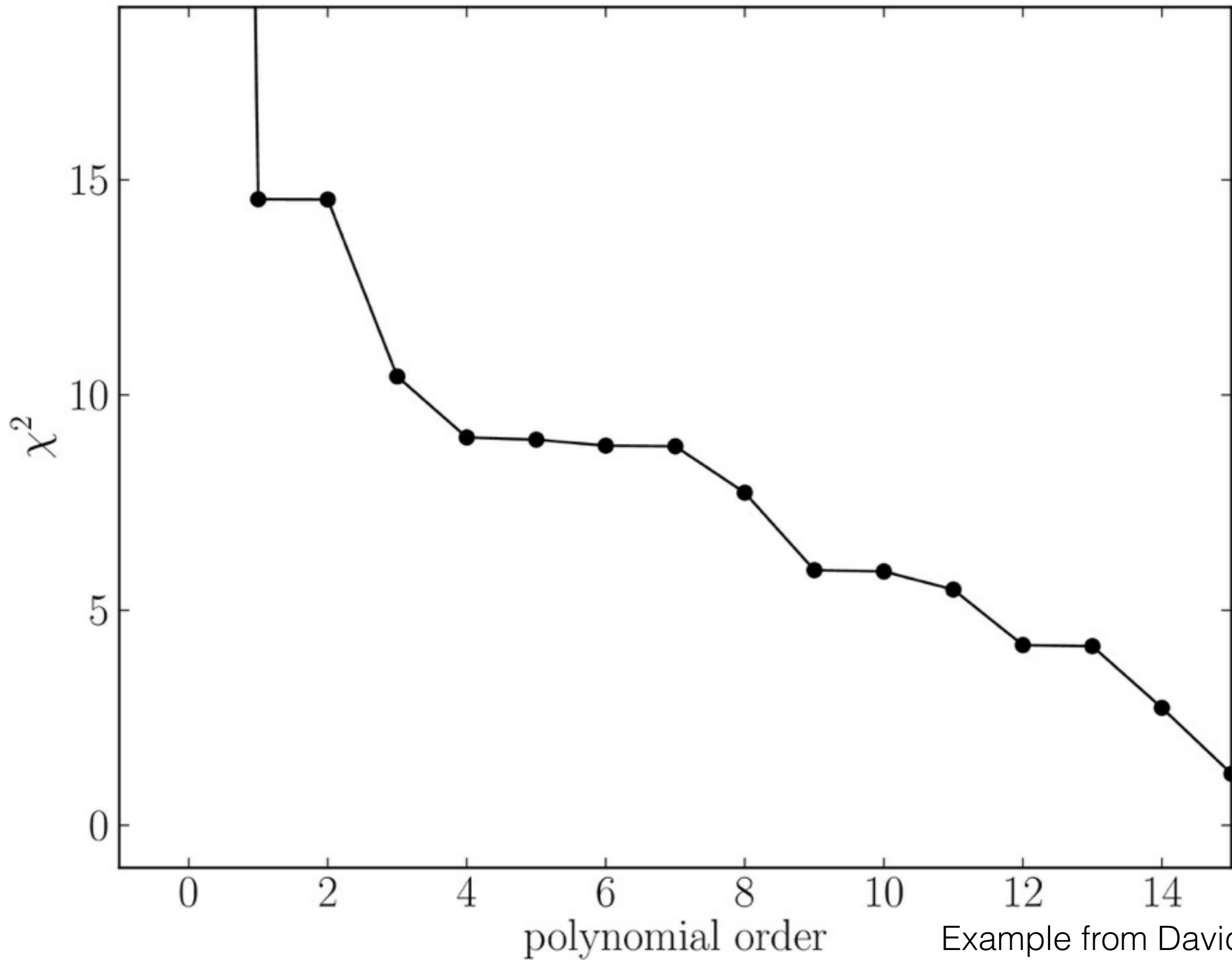
Bayesian Information Criterion (BIC)

- If the posterior is part of the exponential family of PDFs (Gaussian, chi-squared, beta, Bernoulli, ...) can approximate

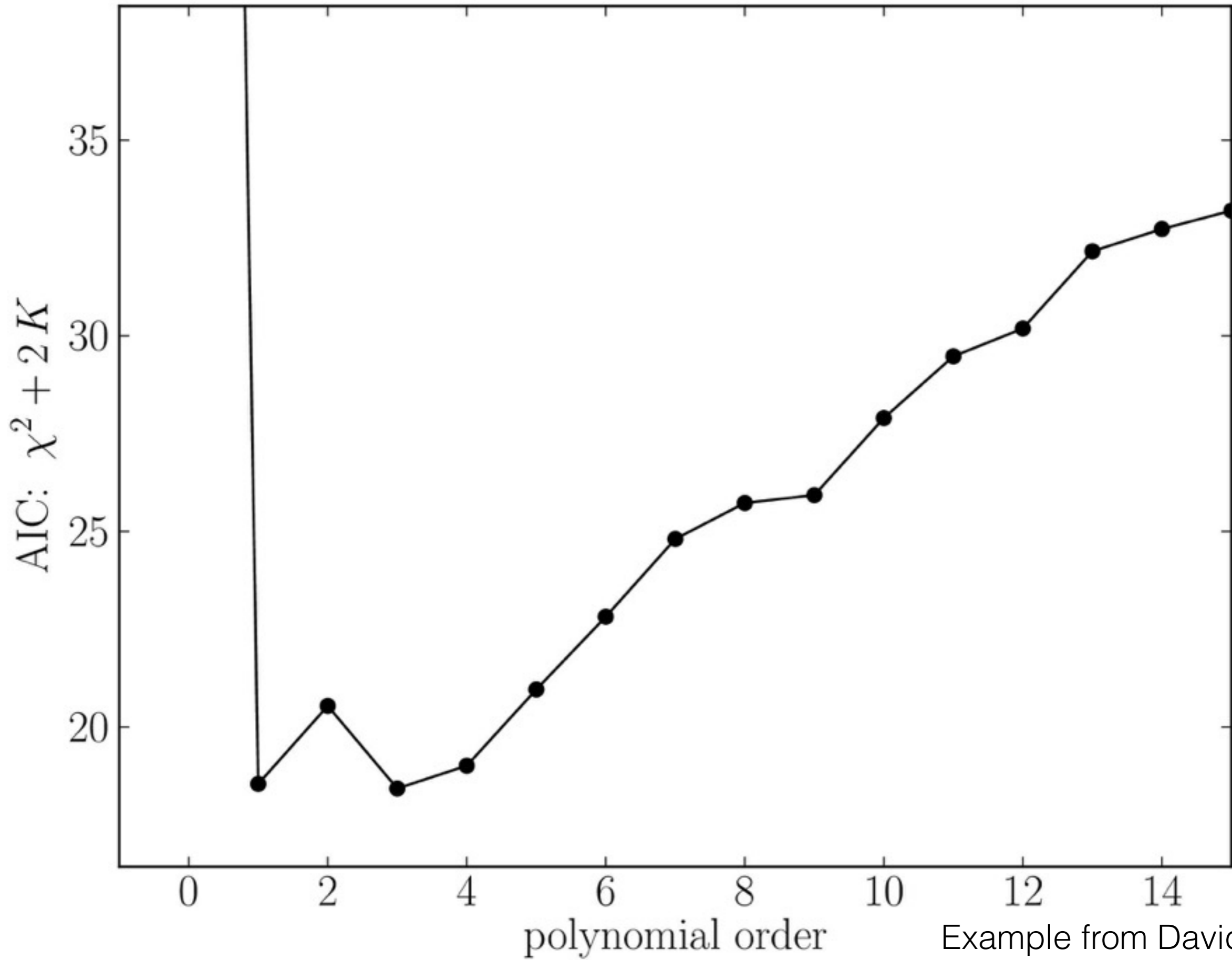
$$\text{Evidence} \sim \text{BIC} = -2 \ln L + K \ln[N]$$

for K parameters and N data points

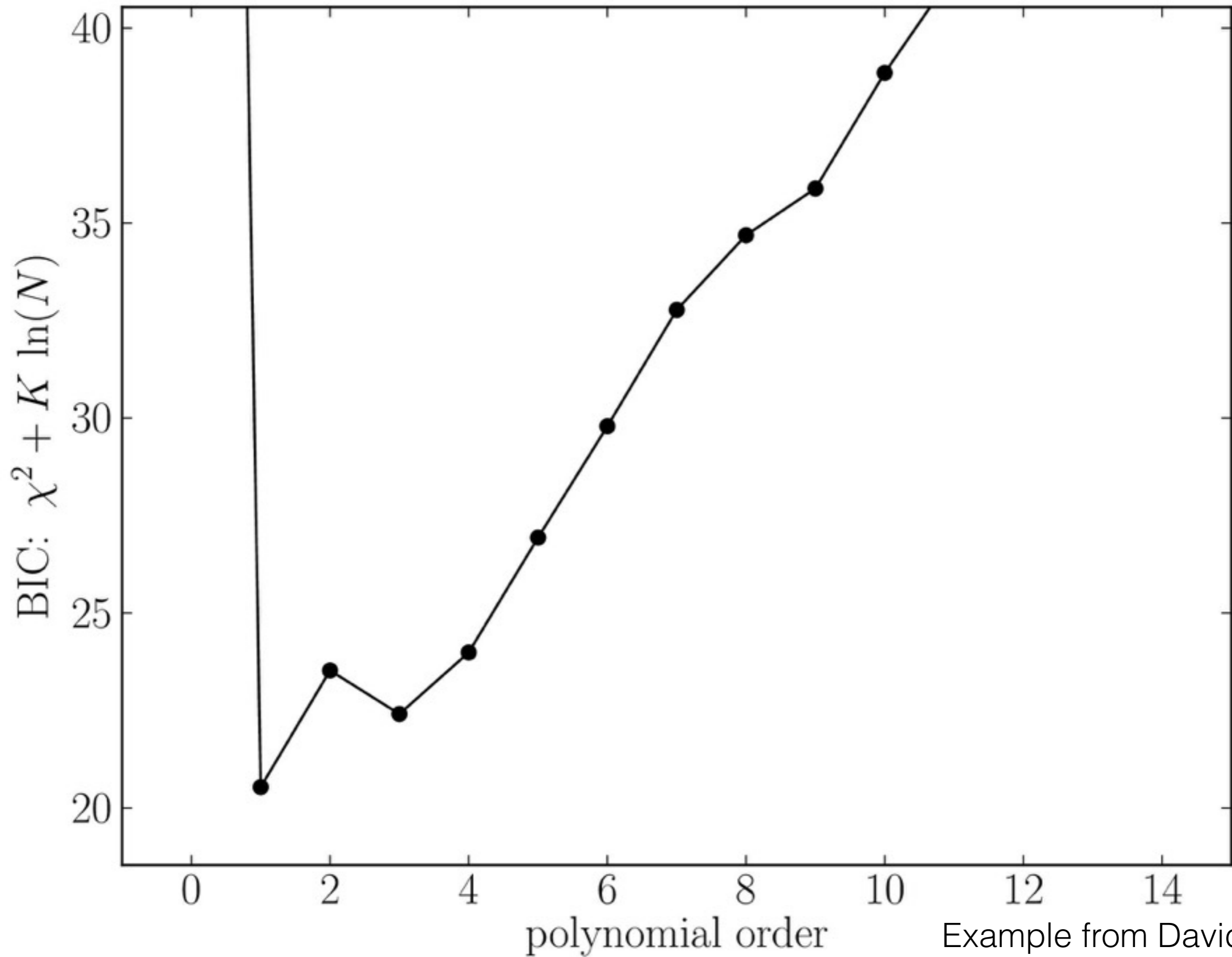
- Similar to AIC, but Bayesian :-)



Example from David Hogg



Example from David Hogg

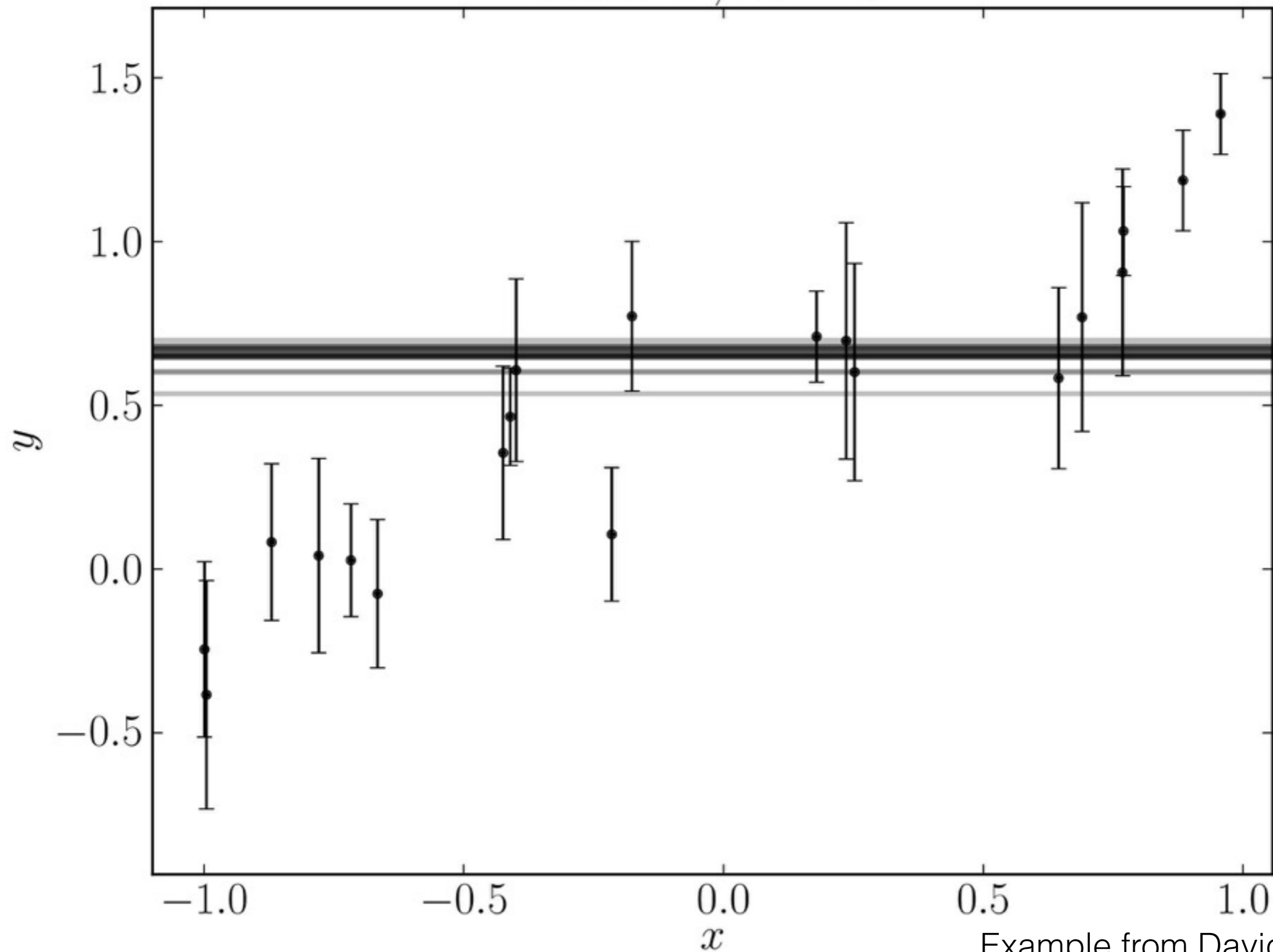


Example from David Hogg

Cross-validation

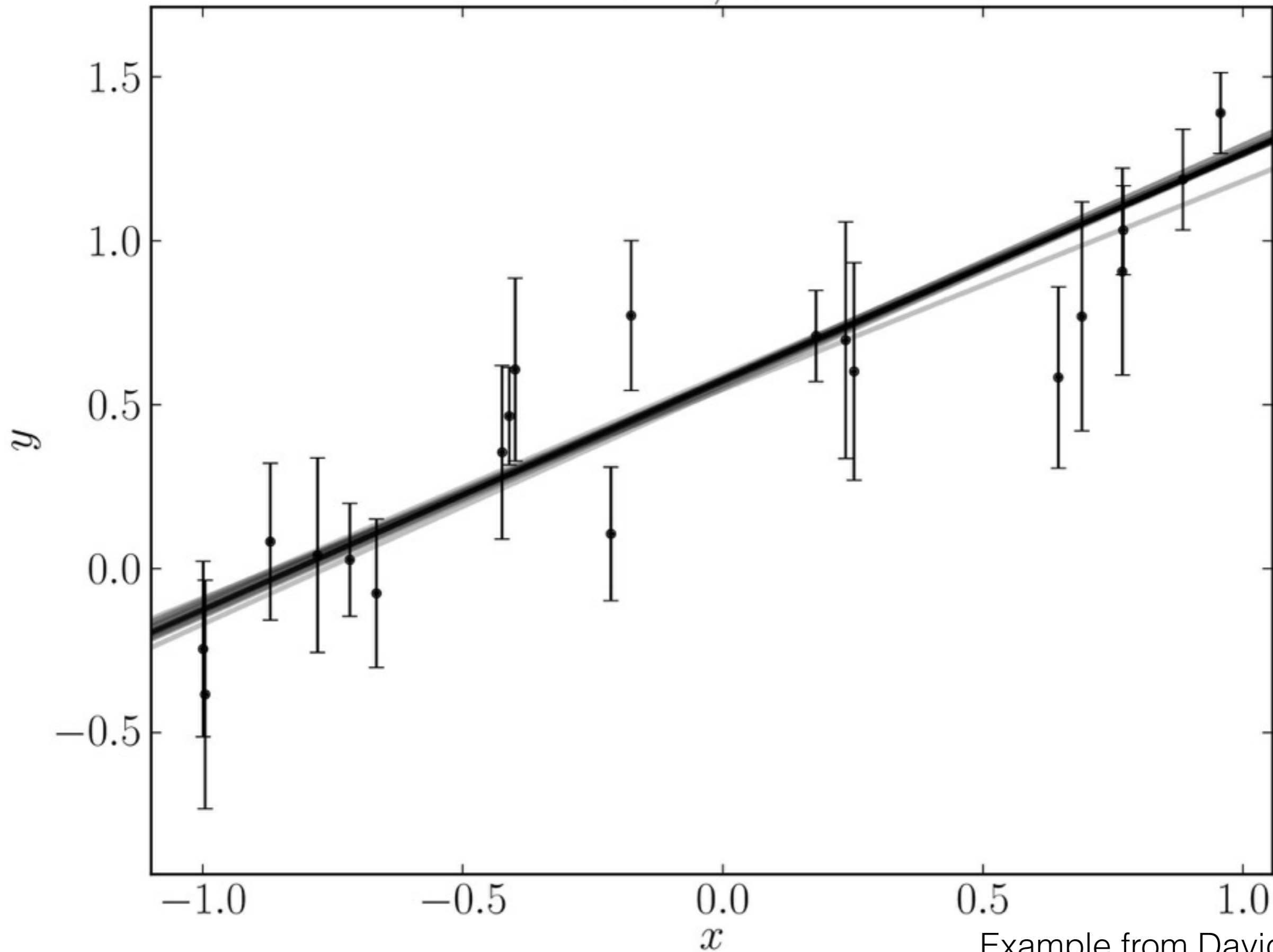
- All previous methods require many assumptions
- Take approach similar to bootstrap and jackknife before, using the data themselves to generate 'new' data
- Cross-validation is very similar to jackknife:
 1. Generate N data sets that leave out 1 data point at a time
[$\{X_1, X_2, X_3, \dots\}$, $\{X_0, X_2, X_3, \dots\}$, $\{X_0, X_1, X_3, \dots\}$, ...]
 2. Fit the model to each data set
 3. Compute the likelihood of the data point that was left out: L_i
 4. Cross validation likelihood $L_{cval} = \text{Prod}_i L_i$

order 0 ; $K = 1$



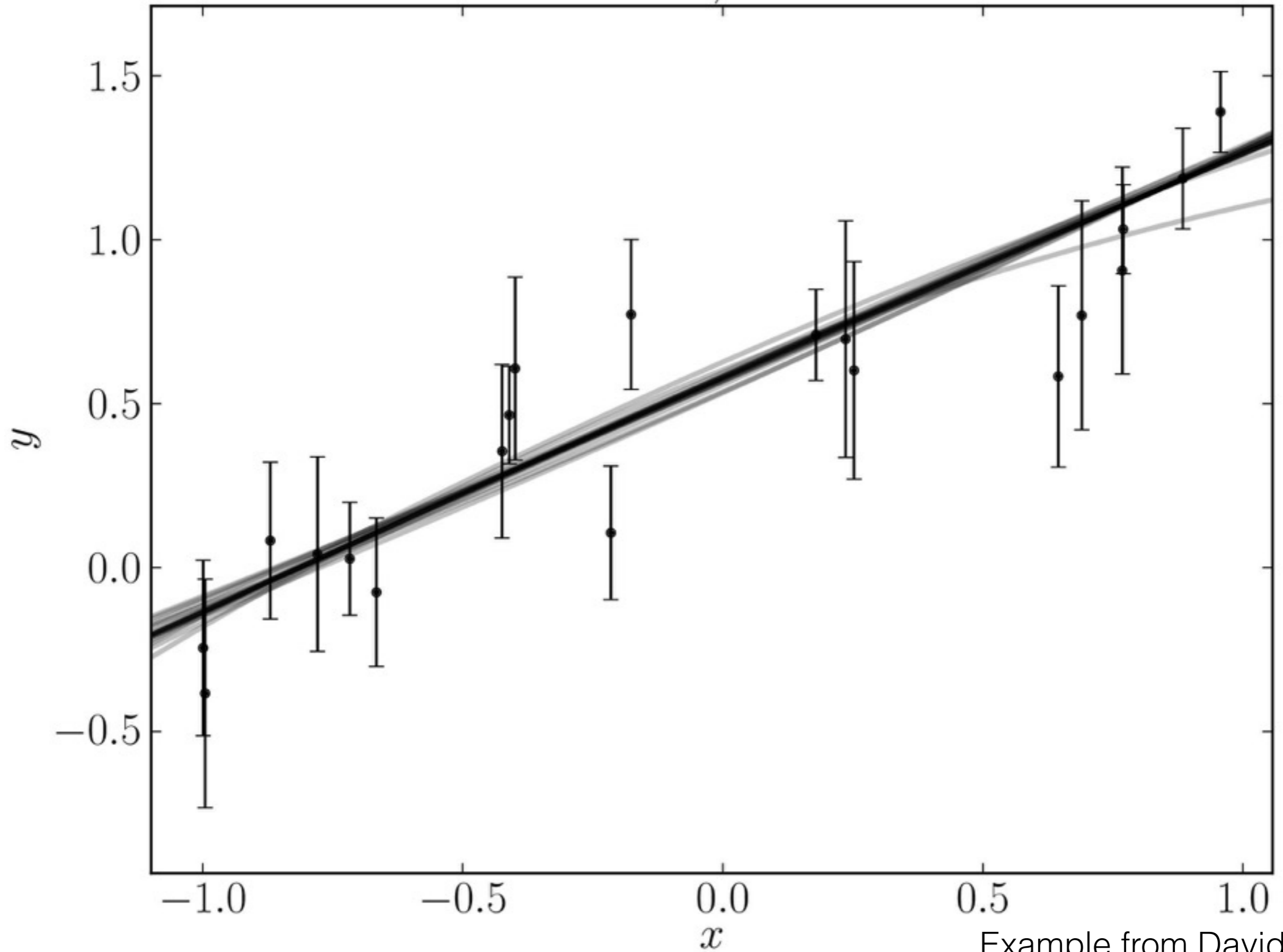
Example from David Hogg

order 1 ; $K = 2$



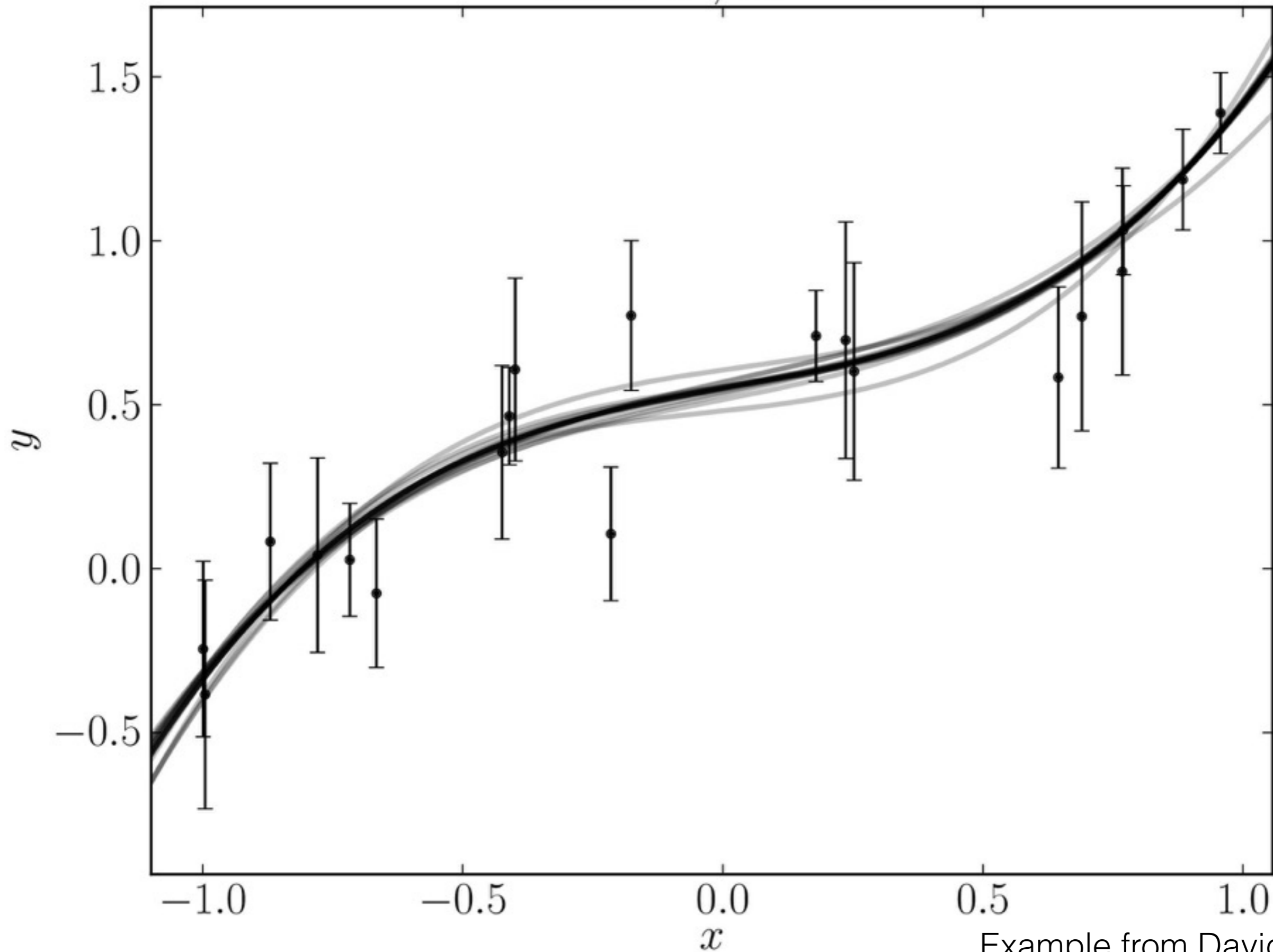
Example from David Hogg

order 2 ; $K = 3$



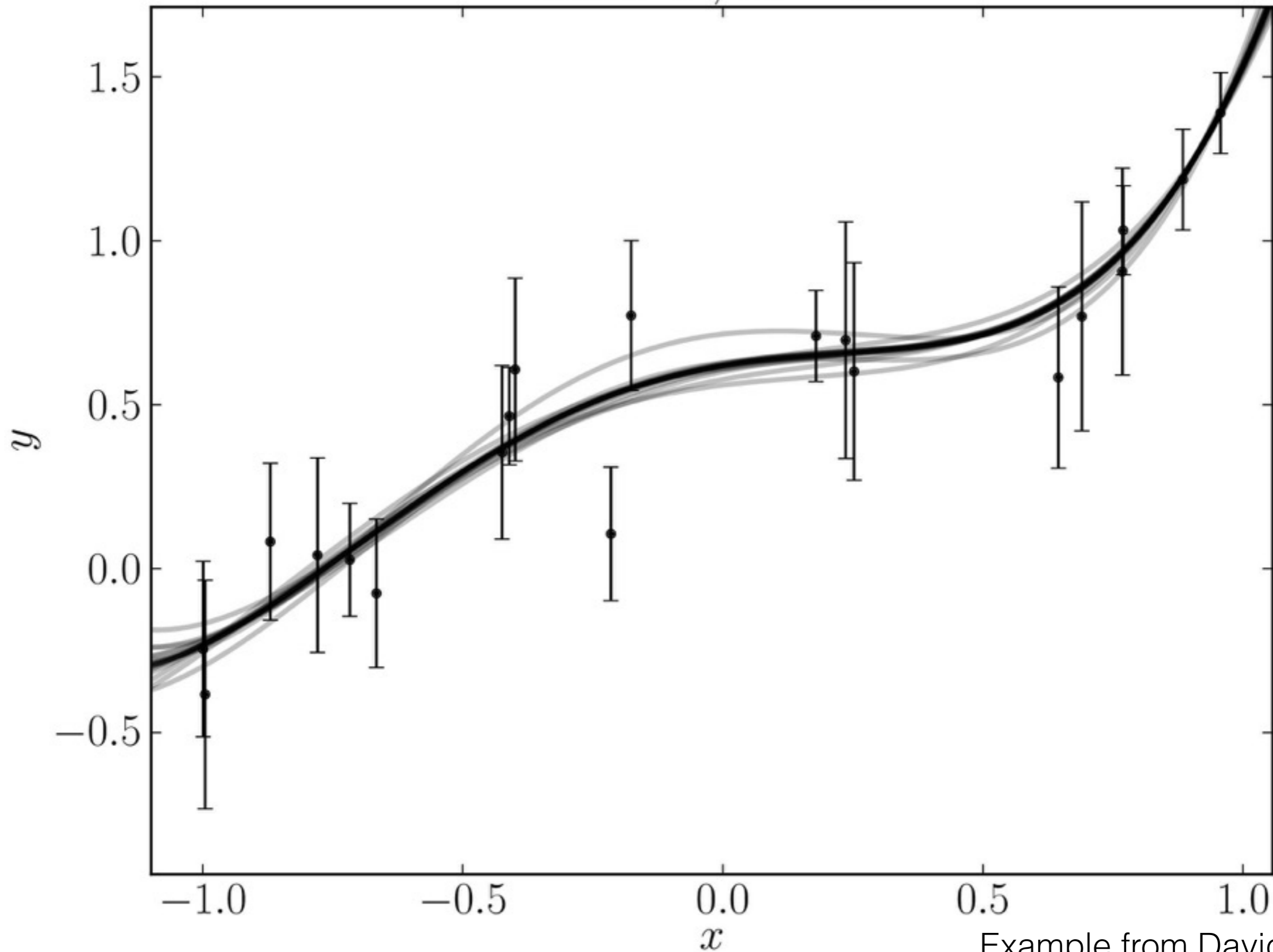
Example from David Hogg

order 3 ; $K = 4$



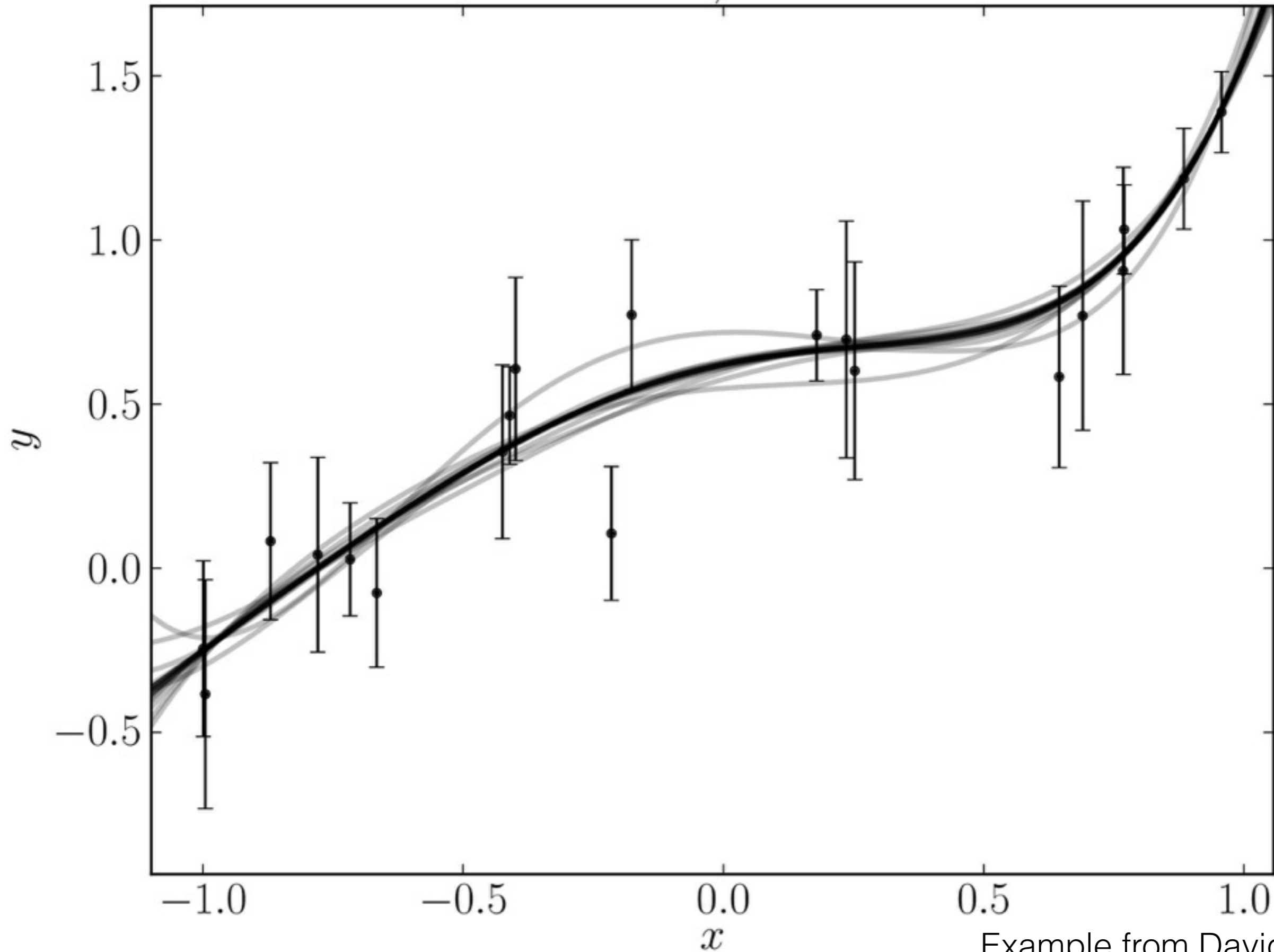
Example from David Hogg

order 4 ; $K = 5$



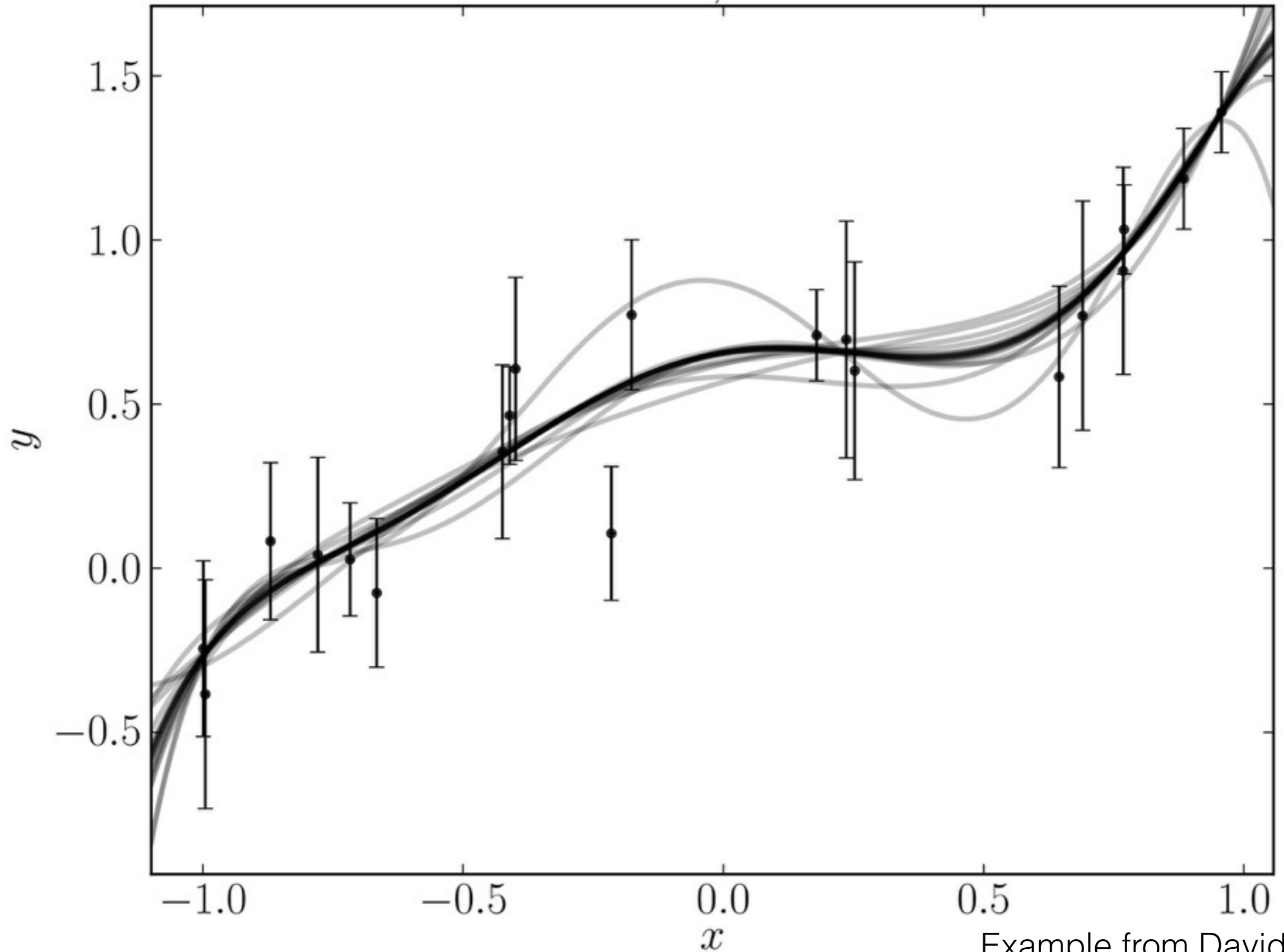
Example from David Hogg

order 5 ; $K = 6$



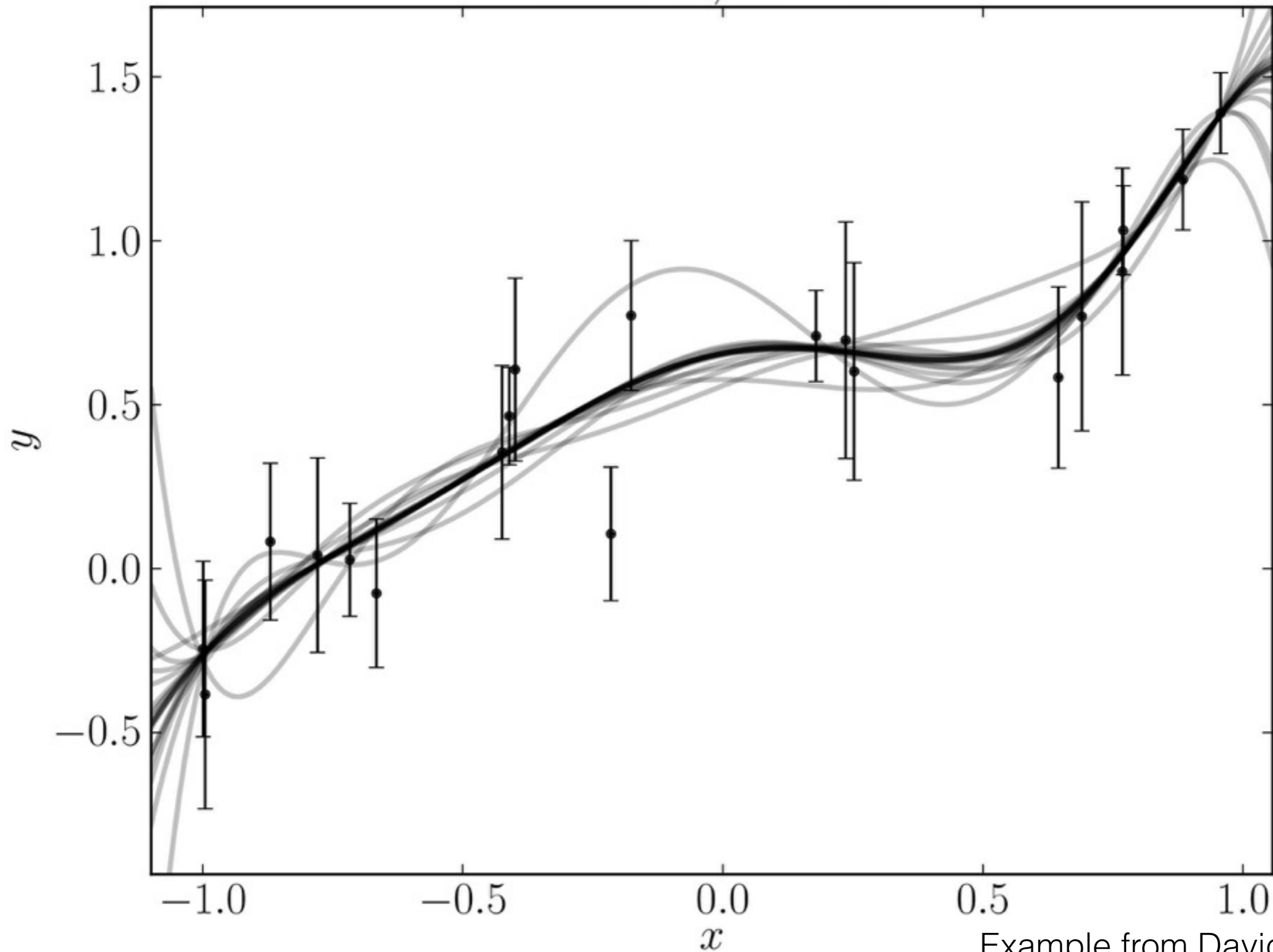
Example from David Hogg

order 6 ; $K = 7$



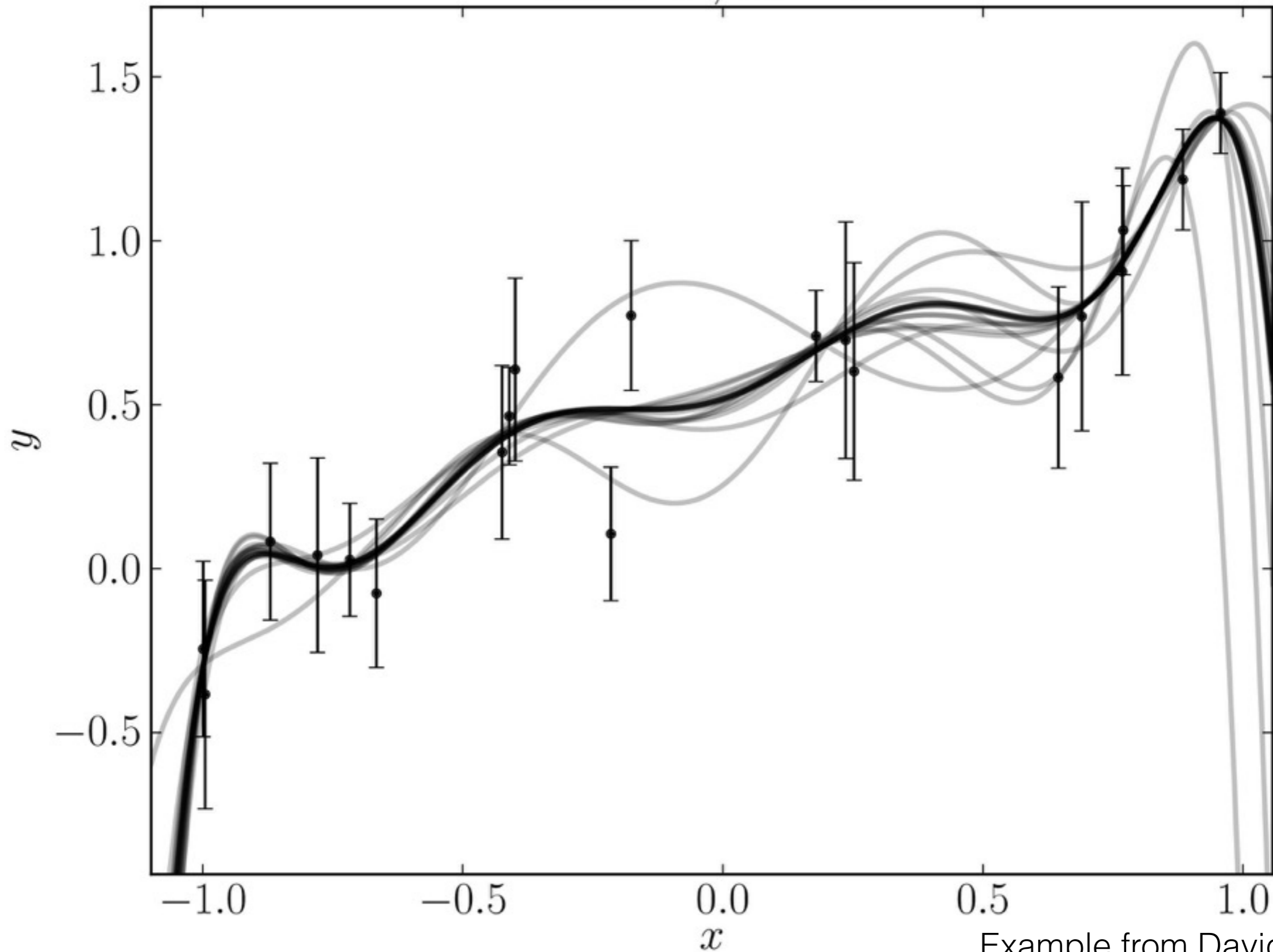
Example from David Hogg

order 7 ; $K = 8$



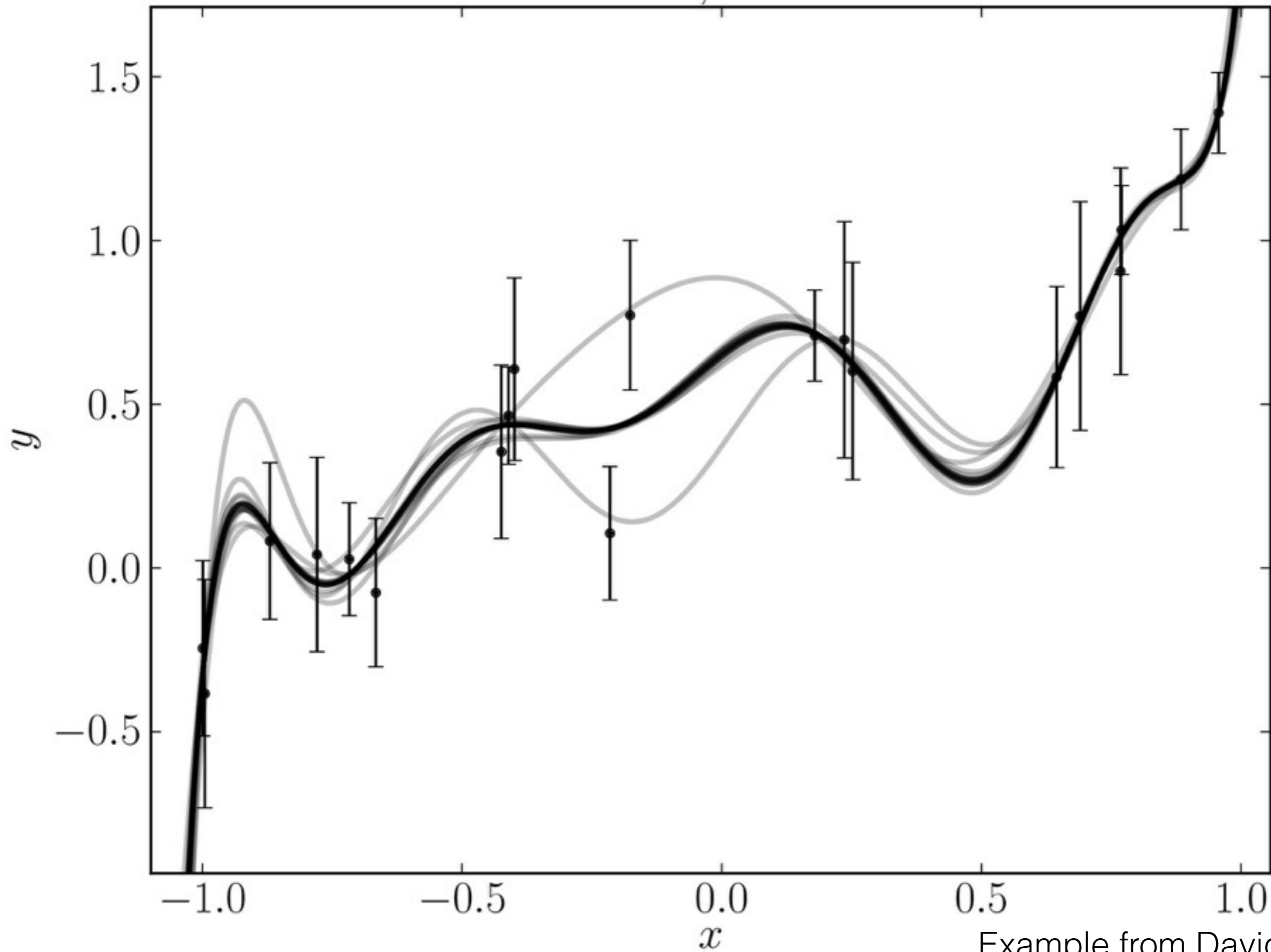
Example from David Hogg

order 8 ; $K = 9$



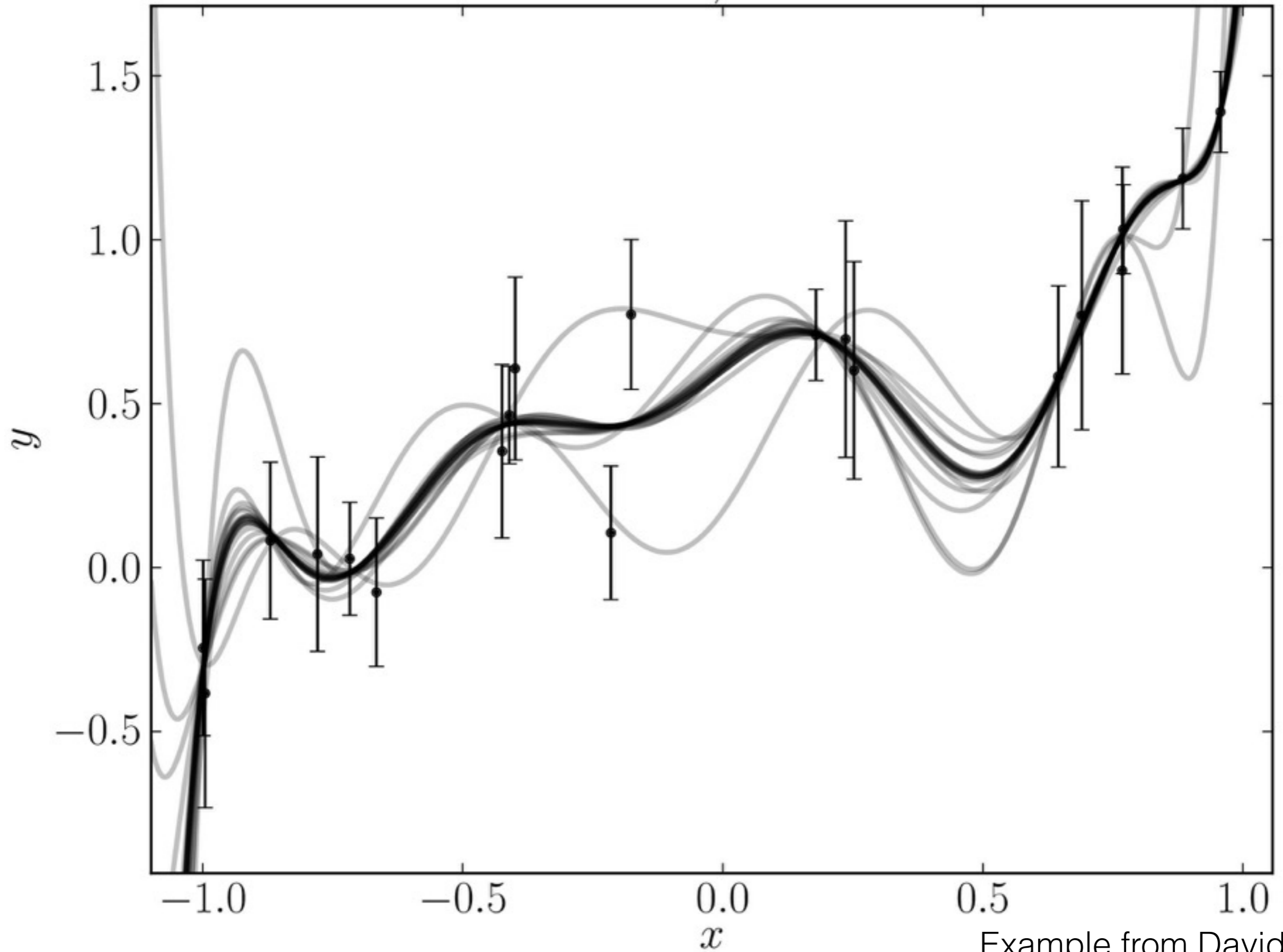
Example from David Hogg

order 9 ; $K = 10$



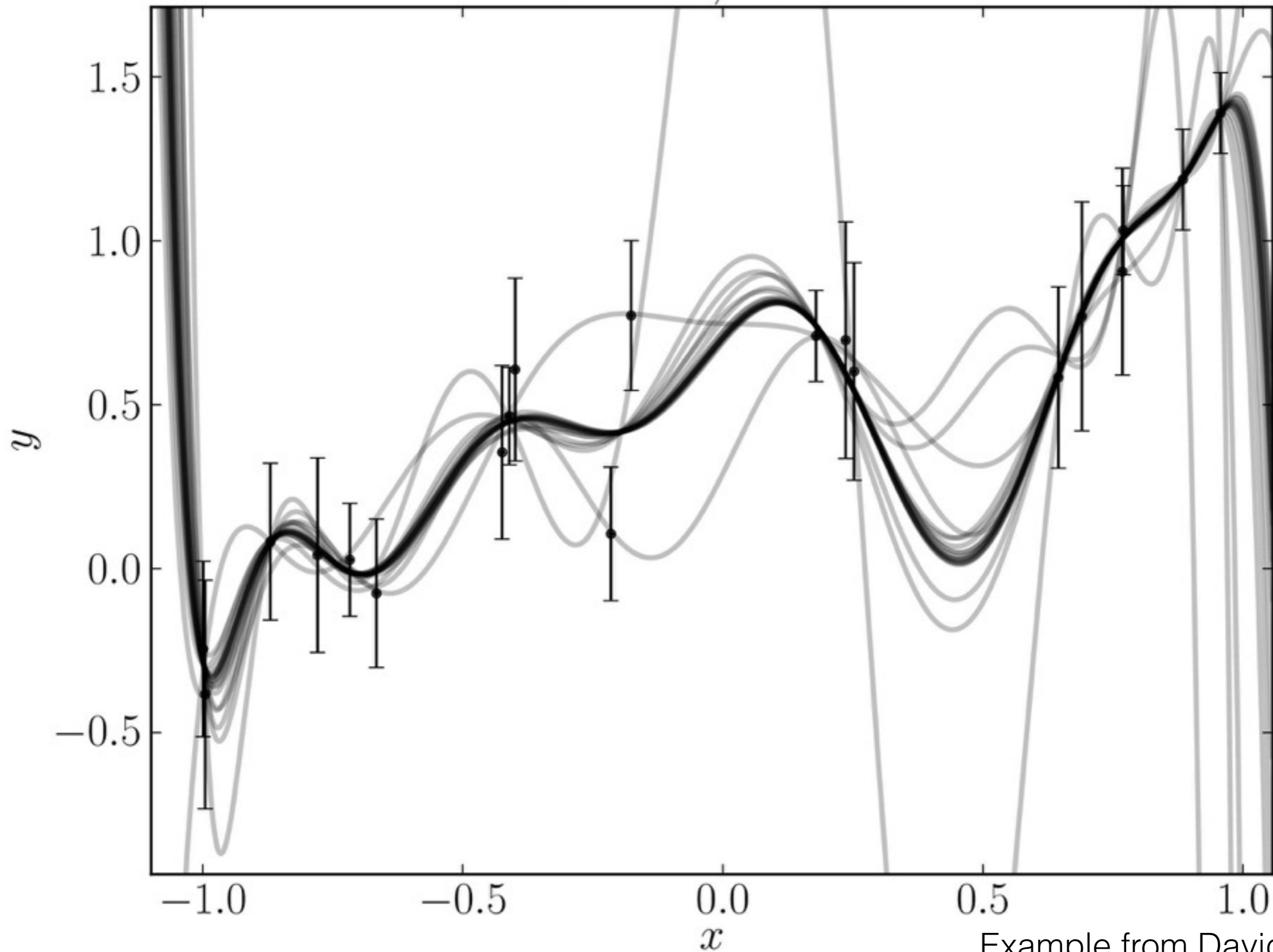
Example from David Hogg

order 10 ; $K = 11$



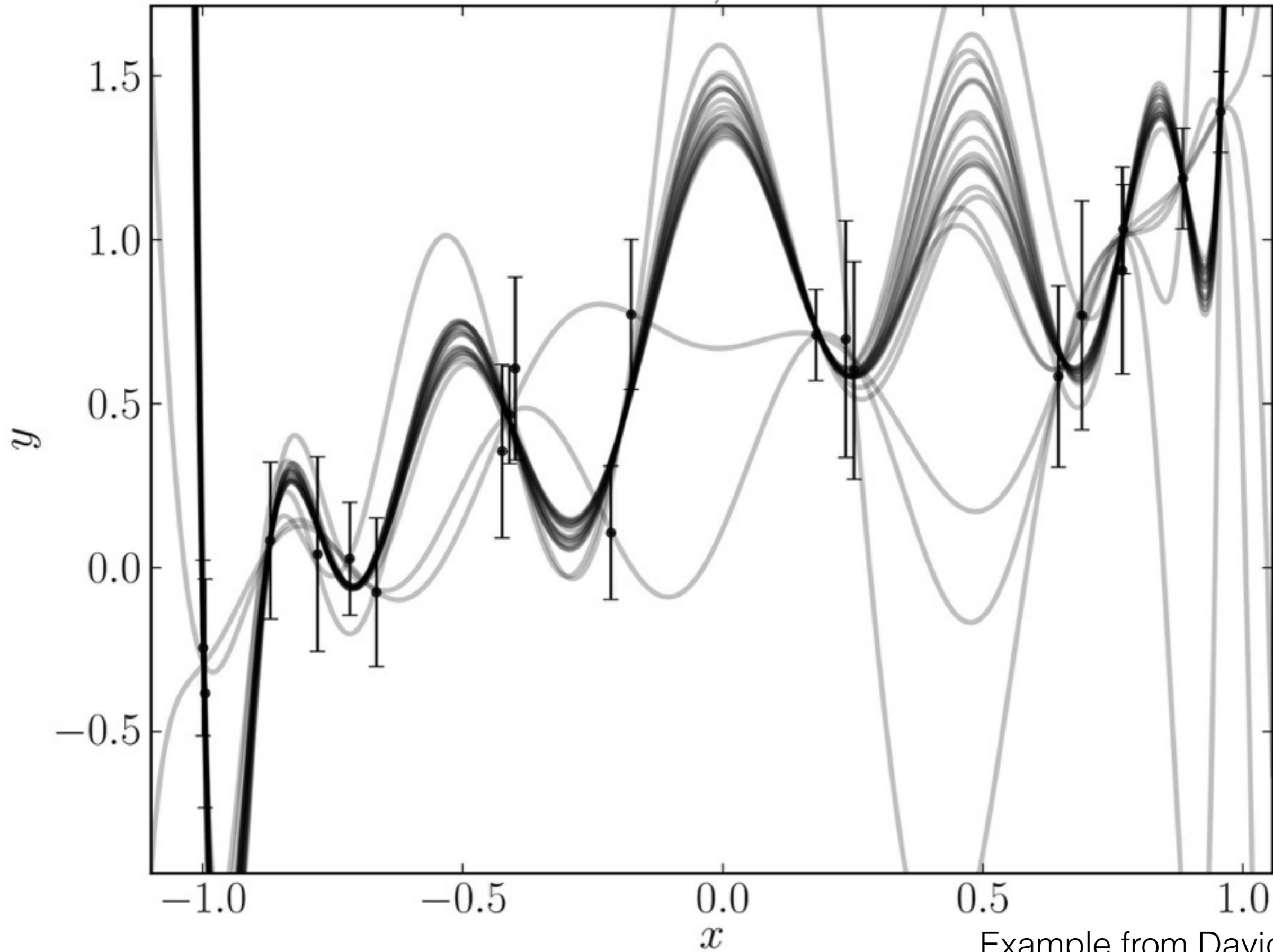
Example from David Hogg

order 11 ; $K = 12$



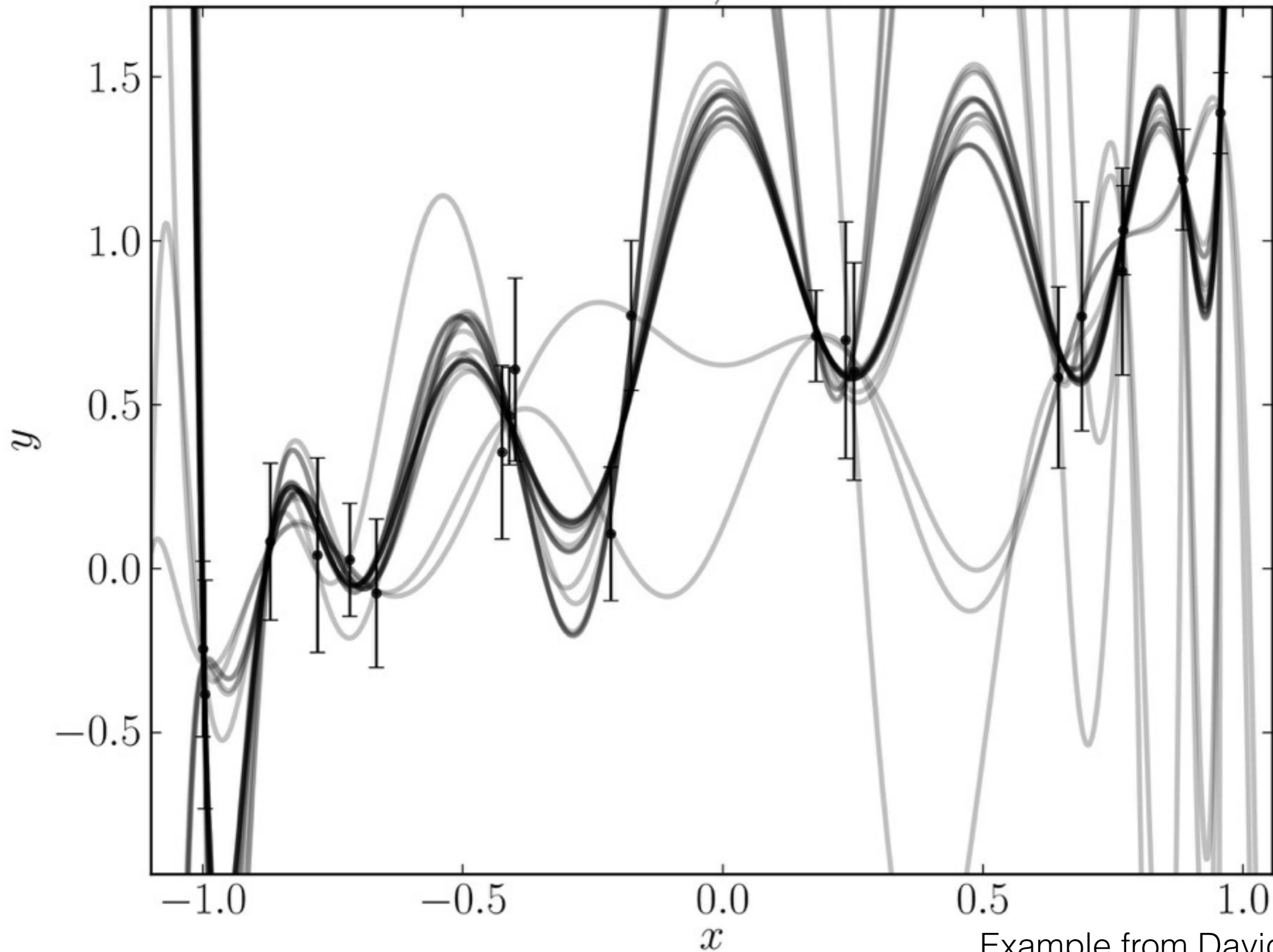
Example from David Hogg

order 12 ; $K = 13$



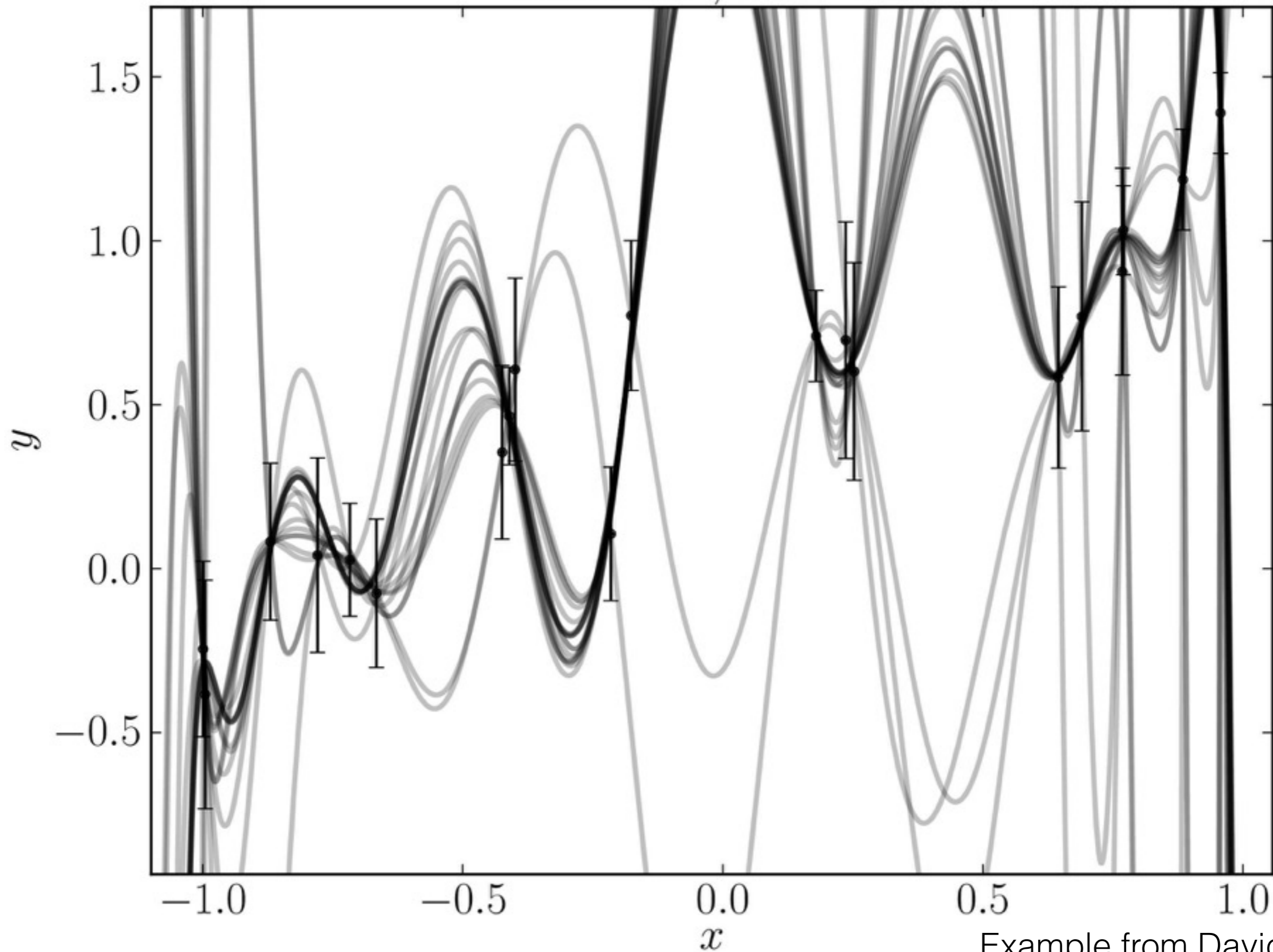
Example from David Hogg

order 13 ; $K = 14$



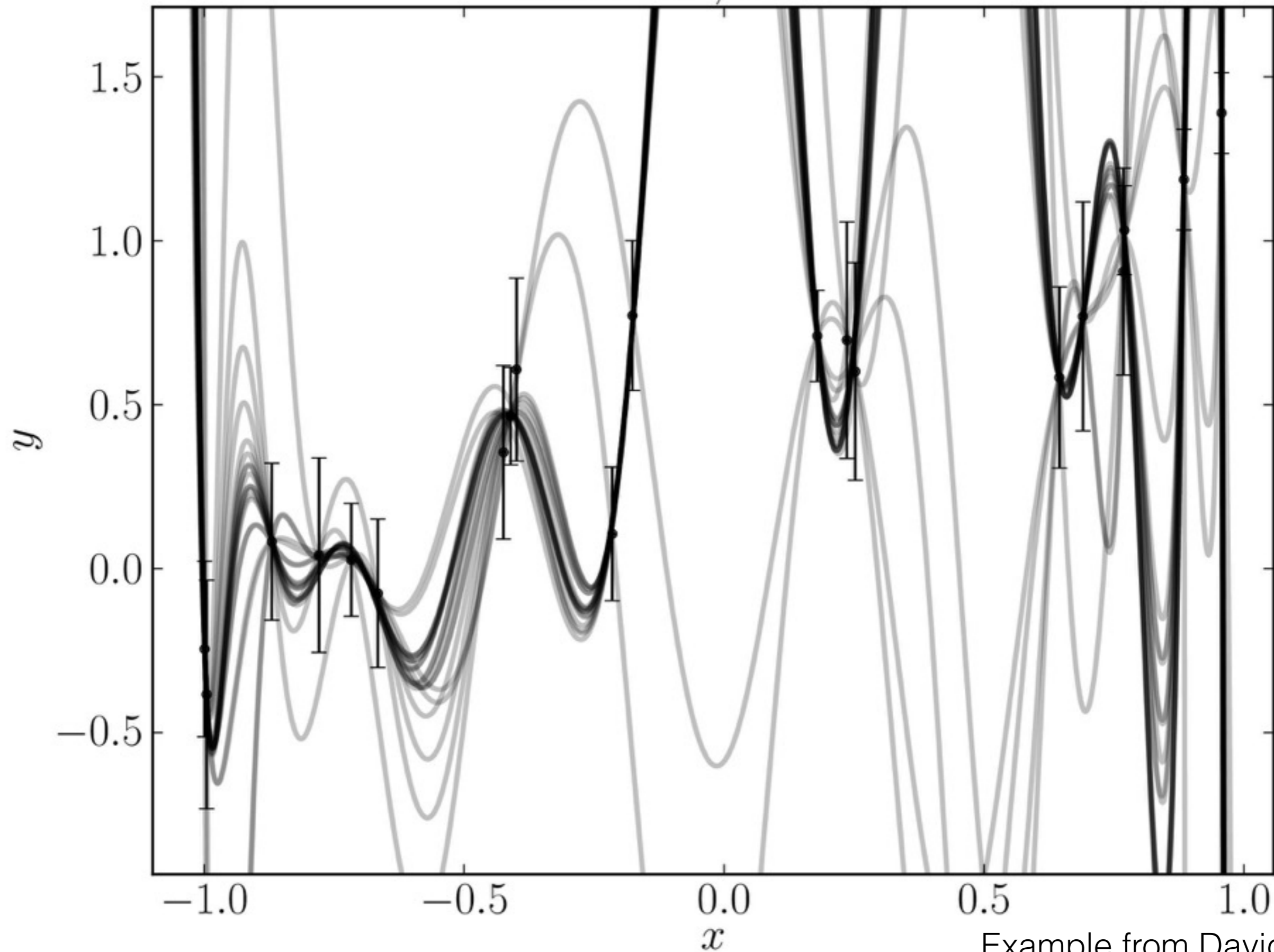
Example from David Hogg

order 14 ; $K = 15$

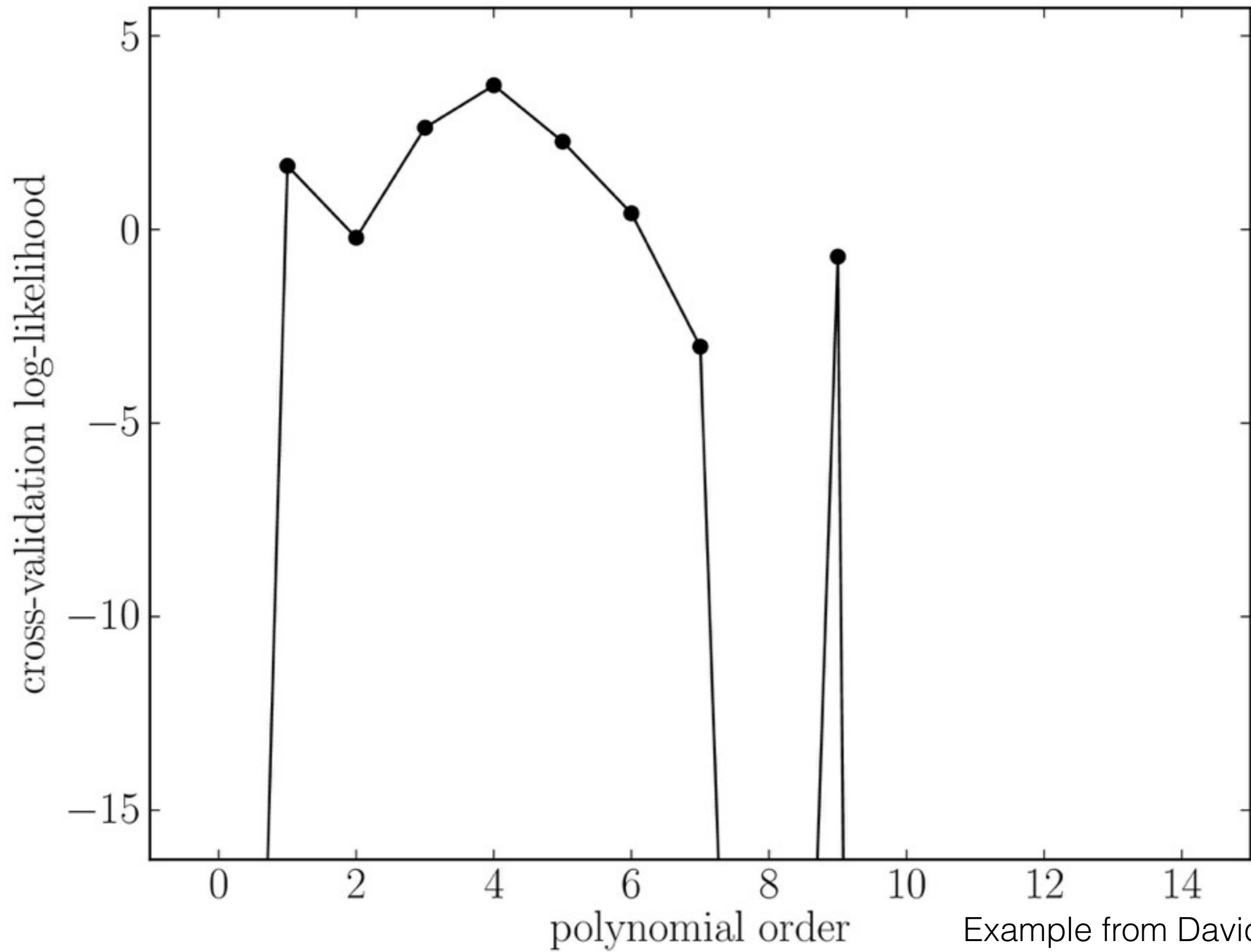


Example from David Hogg

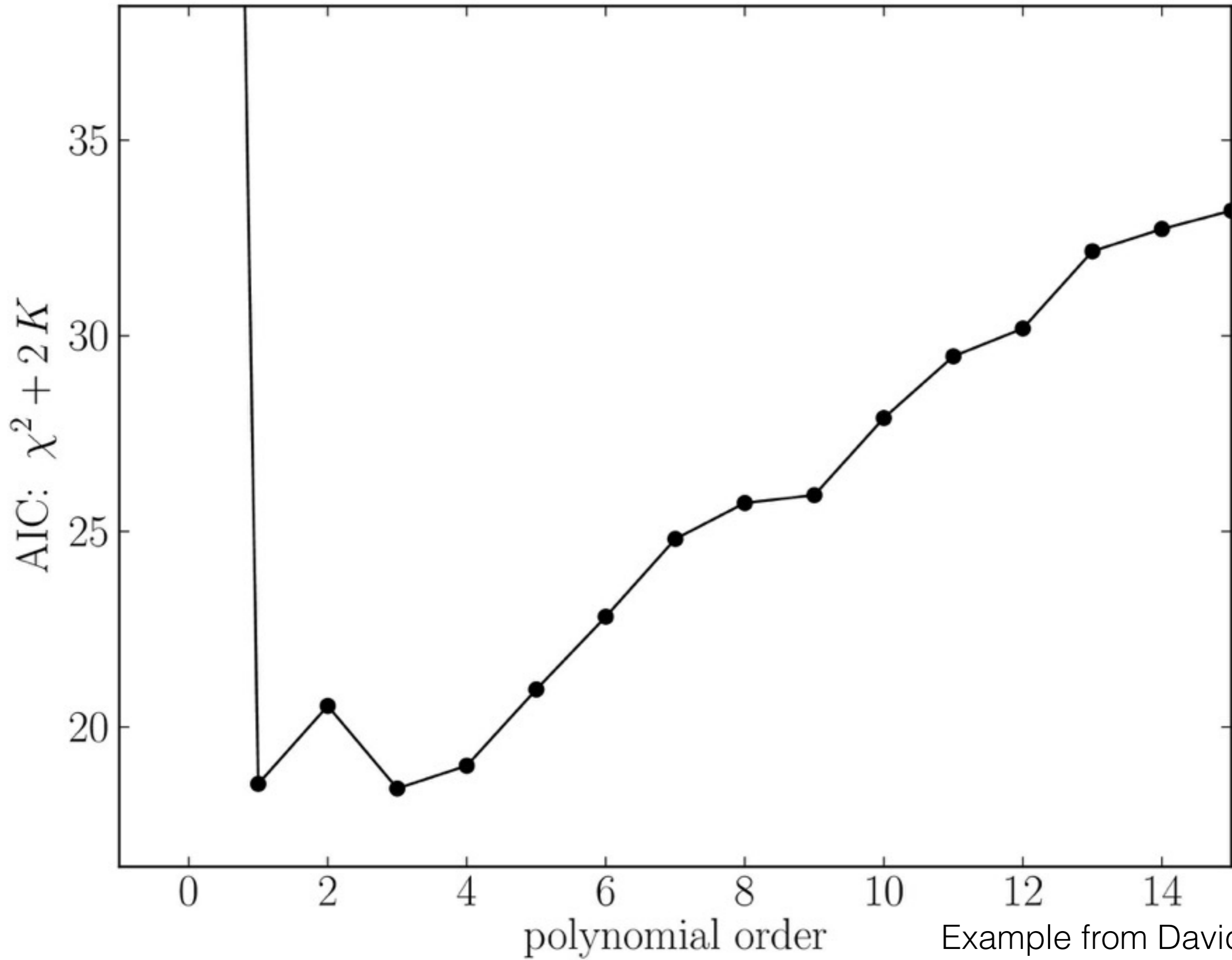
order 15 ; $K = 16$



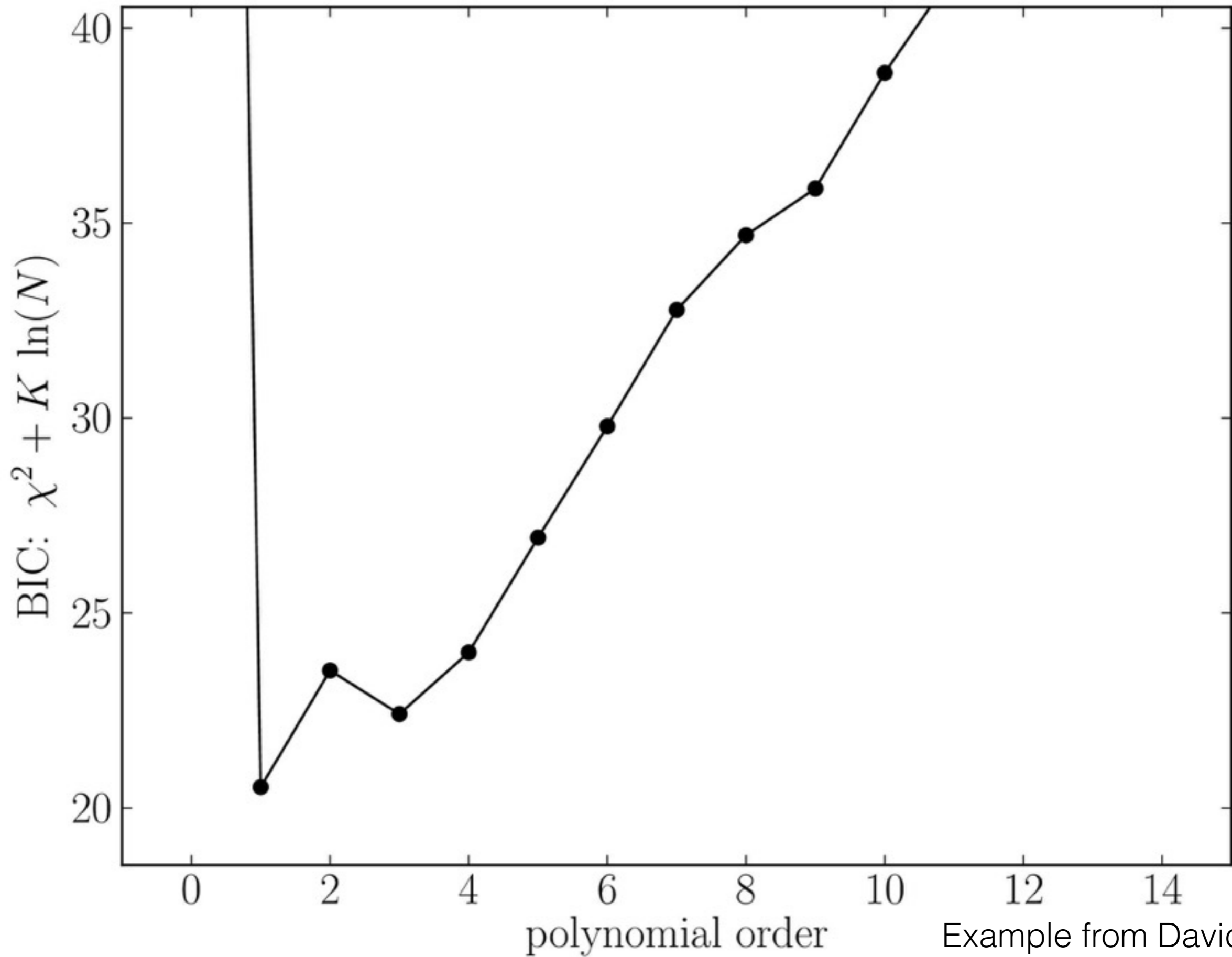
Example from David Hogg



Example from David Hogg

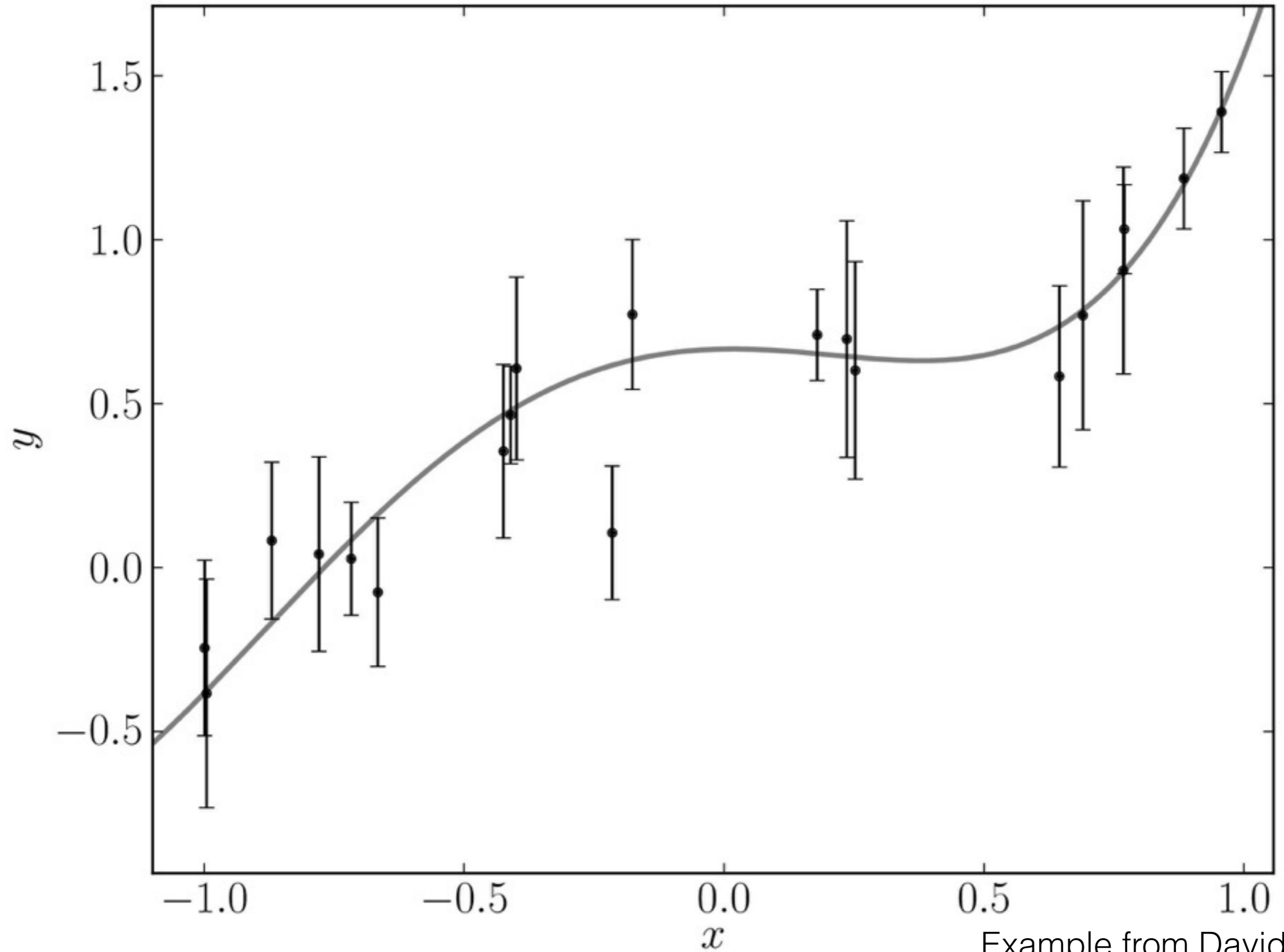


Example from David Hogg



Example from David Hogg

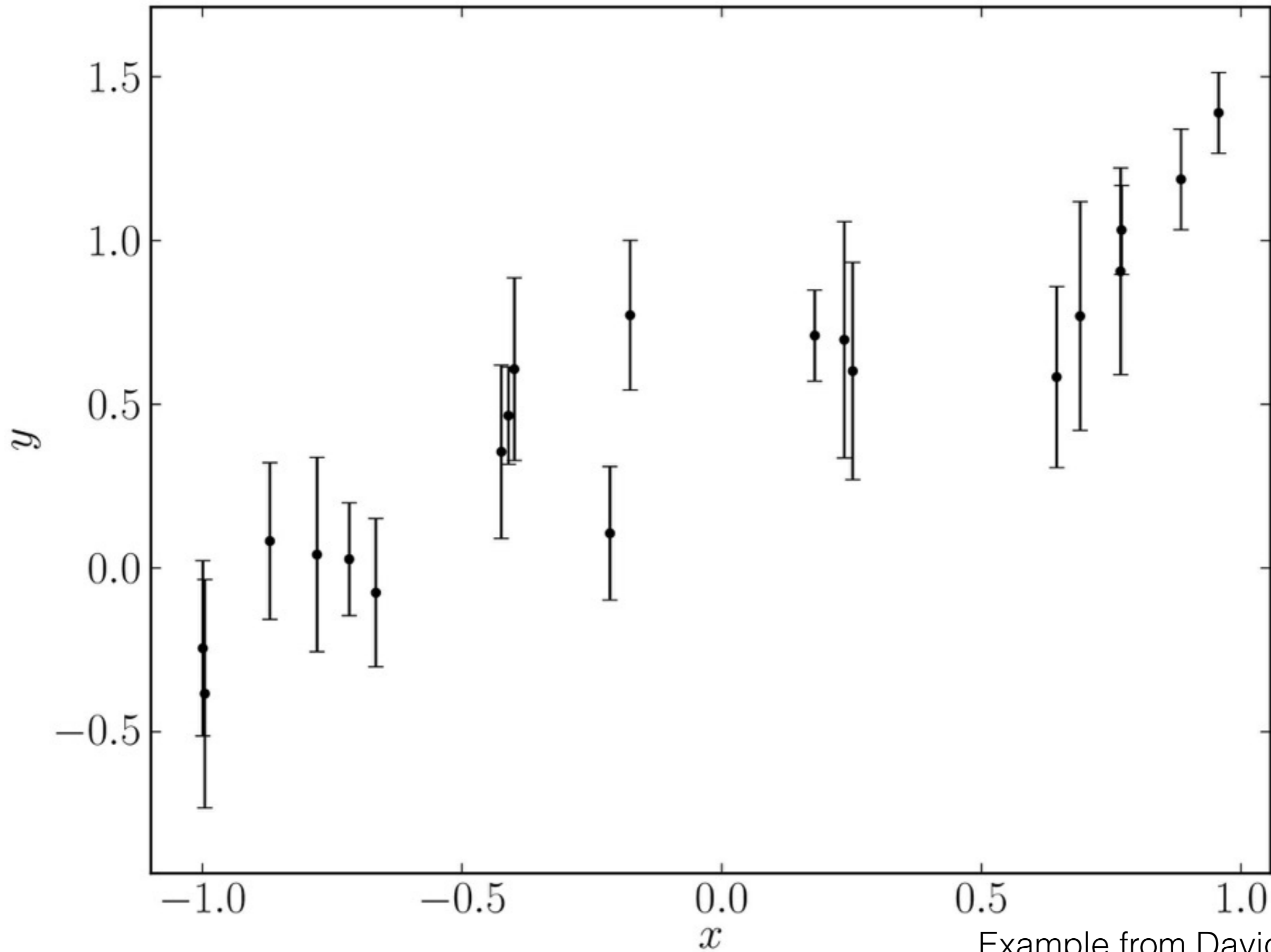
truth: order 4



Example from David Hogg

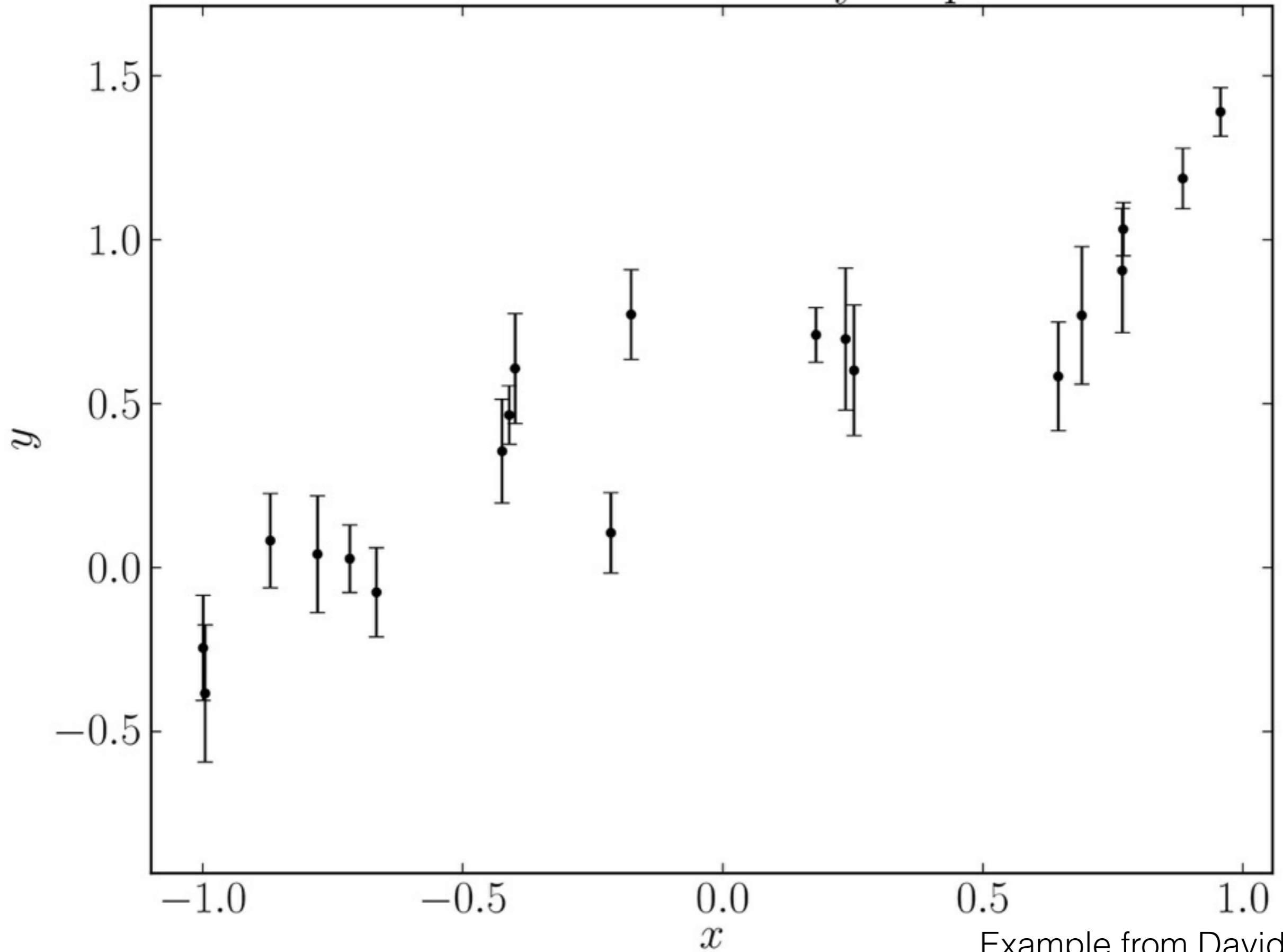
Cross-validation: advantages

- Makes sense and does not make strong assumptions (just that all data points are typical)
- Easy to implement, but can be expensive if each fit is expensive (can use leave-N-out instead)
- Robust against certain types of under/overestimates of the uncertainties (similar to bootstrap)



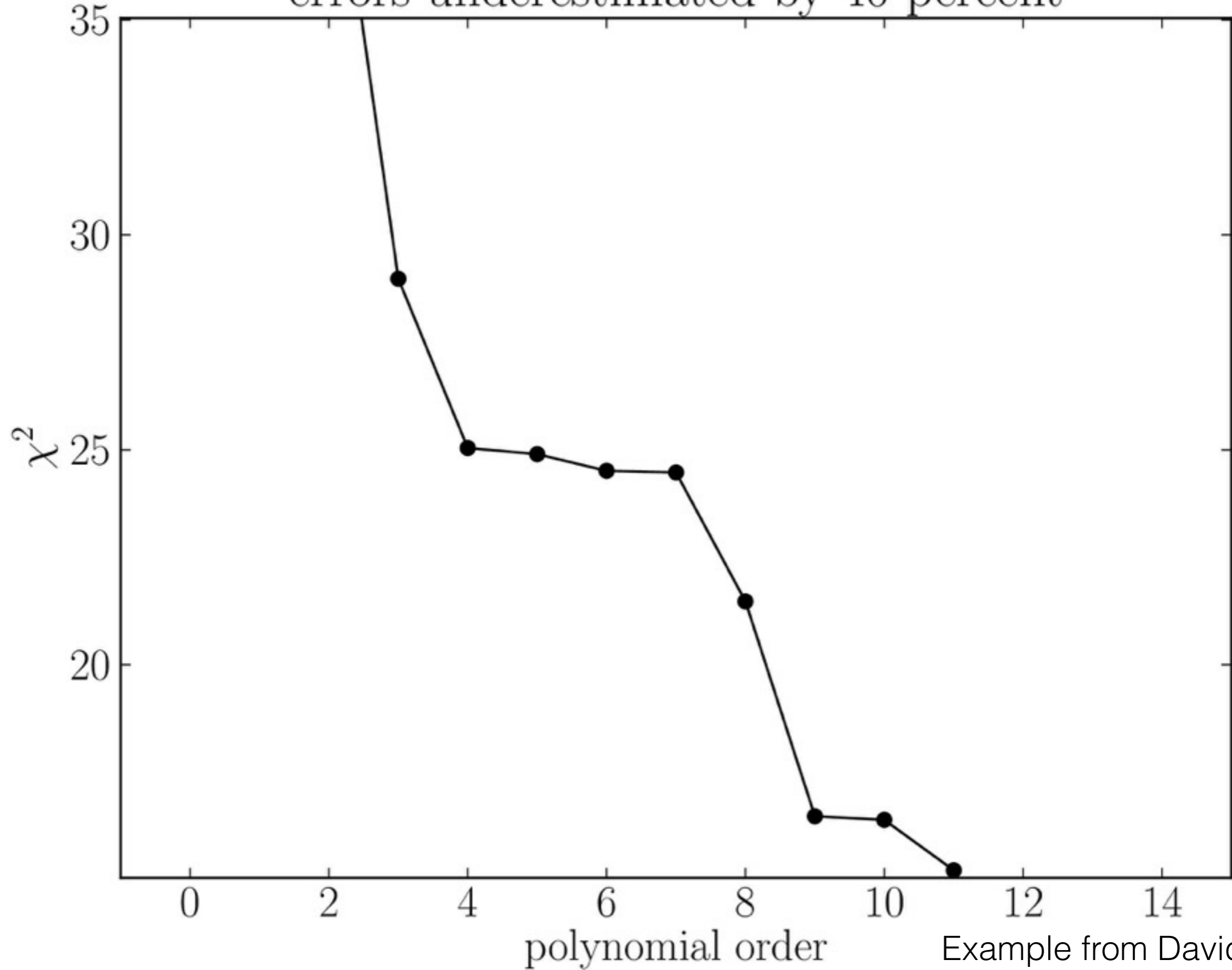
Example from David Hogg

errors underestimated by 40 percent



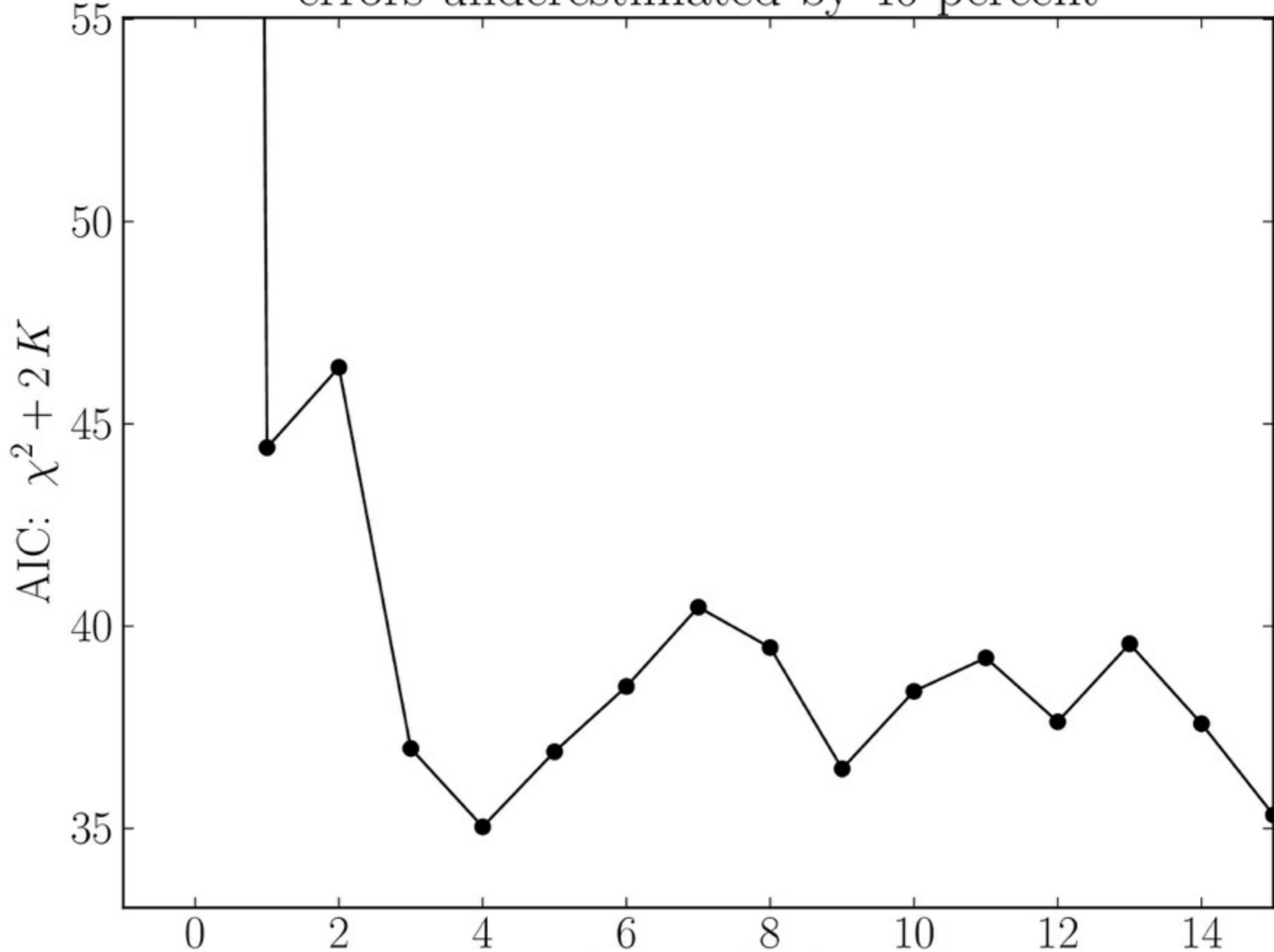
Example from David Hogg

errors underestimated by 40 percent



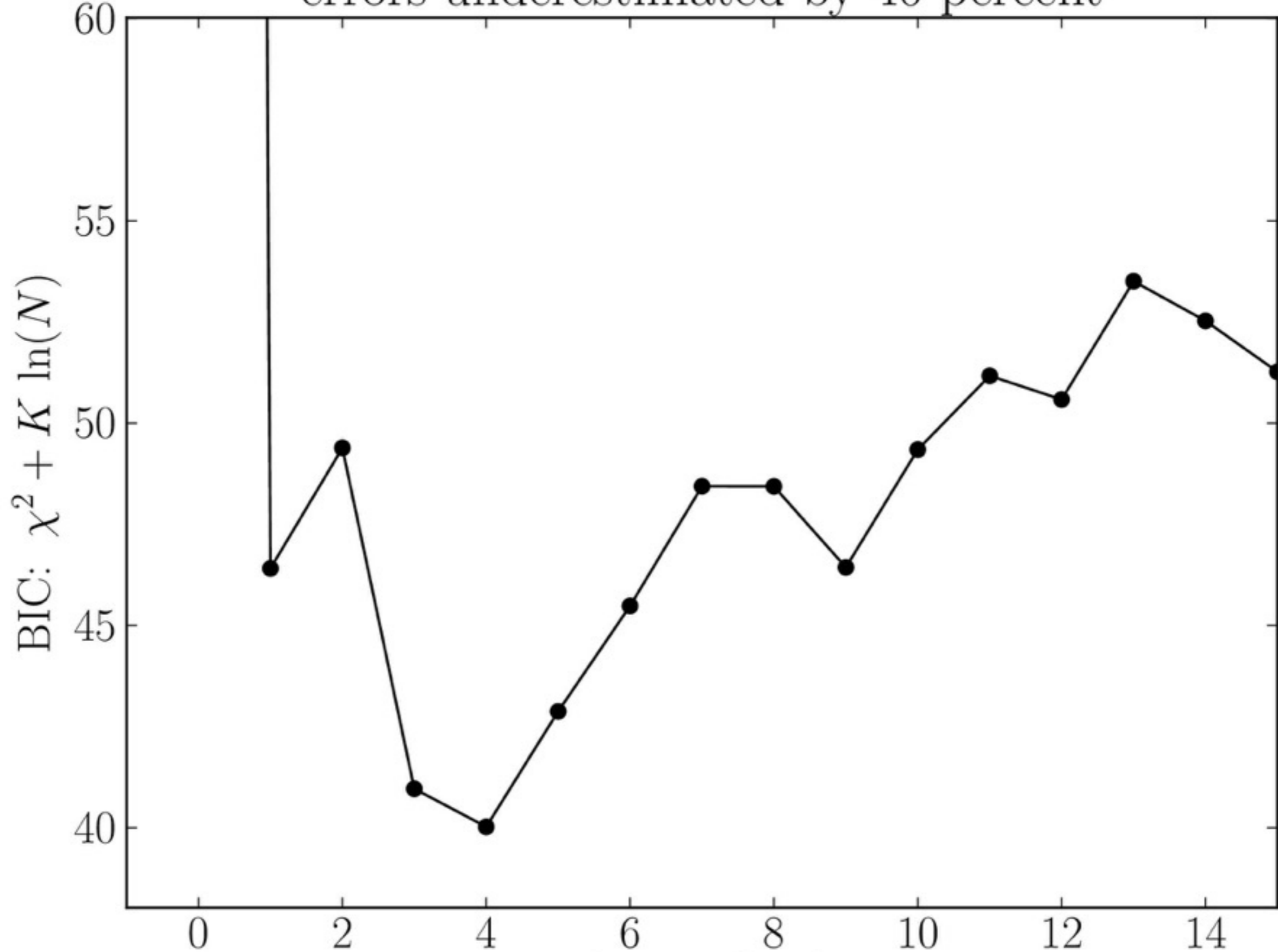
Example from David Hogg

errors underestimated by 40 percent



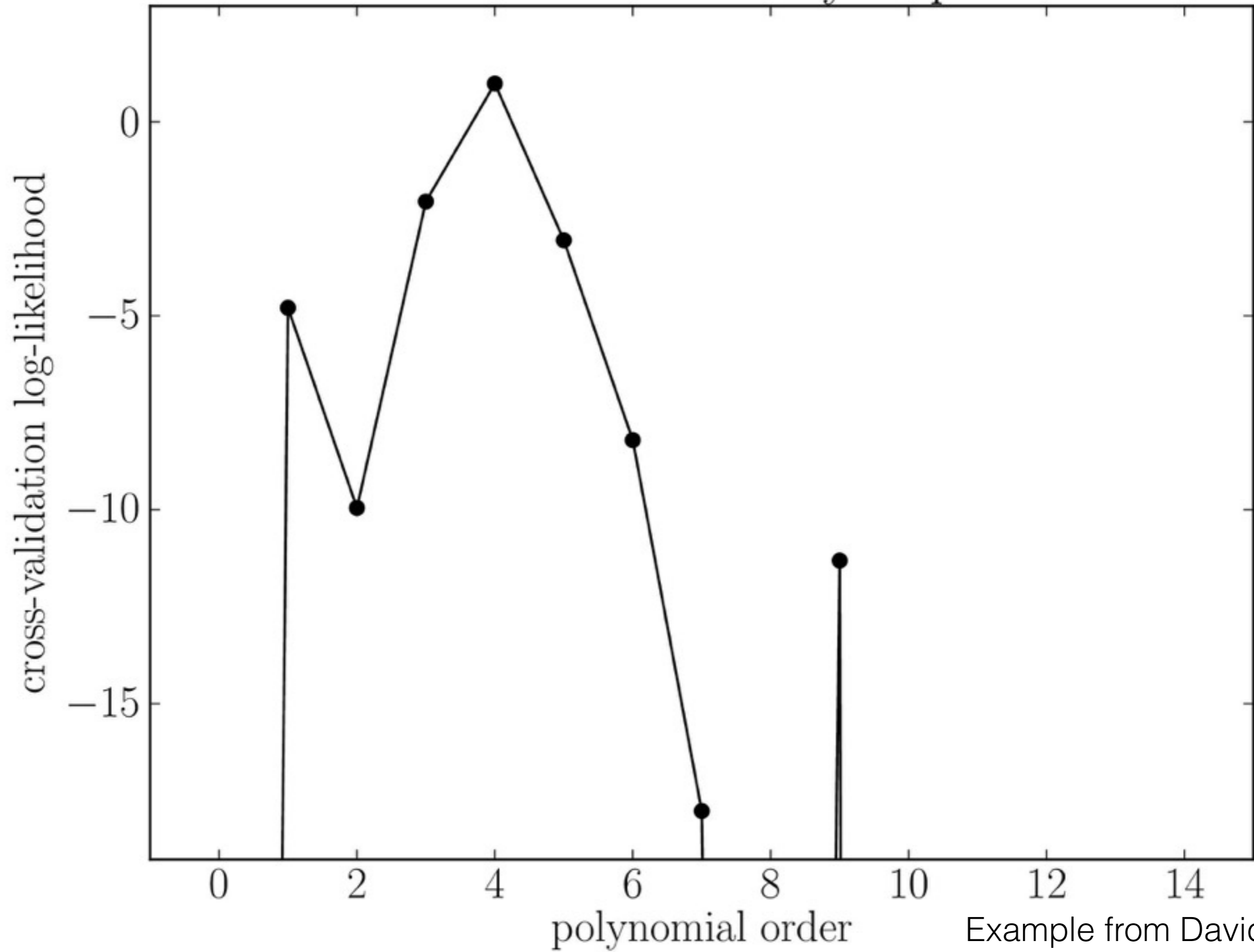
Example from David Hogg

errors underestimated by 40 percent



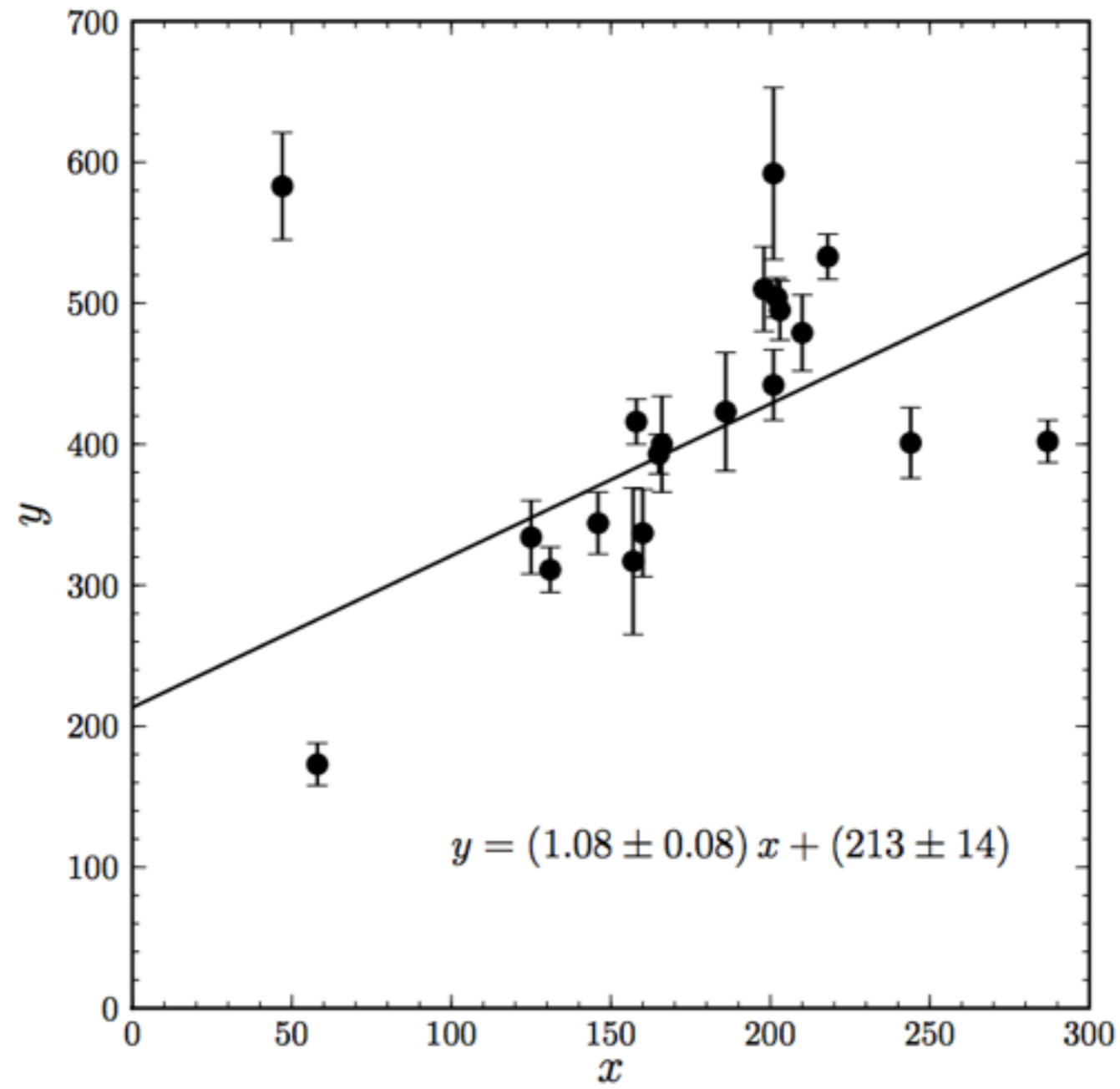
Example from David Hogg

errors underestimated by 40 percent



Example from David Hogg

Outliers, robustics
statist



Handling outliers

- Important! Outliers can wreak havoc to model fitting, model validation, model selection
- Beware of simple cuts. Obvious outliers can be removed (e.g., cosmic rays), but cutting out data points induces selection function near actual data
- Typically best to model outliers w/ simple models
- Or can soften the likelihood for small values (large values of χ^2)

Softened chi²

- Standard chi² strongly punishes large deviations from the mode, can make this smaller
- $X = \sum_i |\text{chi}_i|$
- $\chi_Q^2 = \sum_i \frac{Q^2 [y_i - f(x_i)]^2}{Q^2 \sigma_i^2 + [y_i - f(x_i)]^2}$
- Any heavier-tailed distribution (e.g., student t)
- Changes the likelihood of *every point*, so reduces sensitivity to all points
- Make most sense if you think you do not understand your uncertainties very well

Modeling outliers

- Can use mixture model: some probability that the data are *inliers* and some probability that the data are *outliers*
- Make sense if you think there are two classes of data
- We discussed this in the last class

Mixture model for outliers

In this case, the likelihood is

$$\begin{aligned}\mathcal{L} &\equiv p(\{y_i\}_{i=1}^N | m, b, \{q_i\}_{i=1}^N, Y_b, V_b, I) \\ \mathcal{L} &= \prod_{i=1}^N [p_{\text{fg}}(\{y_i\}_{i=1}^N | m, b, I)]^{q_i} [p_{\text{bg}}(\{y_i\}_{i=1}^N | Y_b, V_b, I)]^{[1-q_i]} \\ \mathcal{L} &= \prod_{i=1}^N \left[\frac{1}{\sqrt{2\pi\sigma_{yi}^2}} \exp\left(-\frac{[y_i - mx_i - b]^2}{2\sigma_{yi}^2}\right) \right]^{q_i} \\ &\quad \times \left[\frac{1}{\sqrt{2\pi[V_b + \sigma_{yi}^2]}} \exp\left(-\frac{[y_i - Y_b]^2}{2[V_b + \sigma_{yi}^2]}\right) \right]^{[1-q_i]},\end{aligned}$$

Posterior requires prior on q_i , introduces new parameter P_b

$$\begin{aligned}p(m, b, \{q_i\}_{i=1}^N, P_b, Y_b, V_b | I) &= p(\{q_i\}_{i=1}^N | P_b, I) p(m, b, P_b, Y_b, V_b | I) \\ p(\{q_i\}_{i=1}^N | P_b, I) &= \prod_{i=1}^N [1 - P_b]^{q_i} P_b^{[1-q_i]},\end{aligned}$$

Mixture model for outliers

- Can be marginalized over individual q_i individually

$$\mathcal{L} \equiv p(\{y_i\}_{i=1}^N | m, b, P_b, Y_b, V_b, I)$$

$$\mathcal{L} \equiv \prod_{i=1}^N [(1 - P_b) p_{\text{fg}}(\{y_i\}_{i=1}^N | m, b, I) + P_b p_{\text{bg}}(\{y_i\}_{i=1}^N | Y_b, V_b, I)]$$

$$\mathcal{L} \propto \prod_{i=1}^N \left[\frac{1 - P_b}{\sqrt{2\pi\sigma_{yi}^2}} \exp\left(-\frac{[y_i - m x_i - b]^2}{2\sigma_{yi}^2}\right) + \frac{P_b}{\sqrt{2\pi[V_b + \sigma_{yi}^2]}} \exp\left(-\frac{[y_i - Y_b]^2}{2[V_b + \sigma_{yi}^2]}\right) \right], \quad (17)$$

Mixture model for outliers

- Can be marginalized over individual q_i individually

$$\mathcal{L} \equiv p(\{y_i\}_{i=1}^N | m, b, P_b, Y_b, V_b, I)$$

$$\mathcal{L} \equiv \prod_{i=1}^N [(1 - P_b) p_{\text{fg}}(\{y_i\}_{i=1}^N | m, b, I) + P_b p_{\text{bg}}(\{y_i\}_{i=1}^N | Y_b, V_b, I)]$$

$$\mathcal{L} \propto \prod_{i=1}^N \left[\frac{1 - P_b}{\sqrt{2\pi\sigma_{yi}^2}} \exp\left(-\frac{[y_i - m x_i - b]^2}{2\sigma_{yi}^2}\right) + \frac{P_b}{\sqrt{2\pi[V_b + \sigma_{yi}^2]}} \exp\left(-\frac{[y_i - Y_b]^2}{2[V_b + \sigma_{yi}^2]}\right) \right] \quad (17)$$

Inliers

Outliers

Mixture model for outliers

- Works well in general:
 - Simple model, e.g., Gaussian with some free mean and variance
 - Halo when looking at disk kinematics
 - ...
- More expensive than softening the objective function, because additional parameters, but retain more information

Robust statistics

- Most frequentist and Bayesian methods have a lot of assumptions about the distribution of the data
- Robust statistic is robust to deviations in the assumptions, e.g., deviations from Gaussian uncertainties, mis-estimated uncertainties, outliers
- Typically better to properly model one's uncertainties and outliers, but robust statistics useful in general data handling

Median

- Mean is a bad estimator for the central location of the data if the uncertainties are non-Gaussian and have long-tail (~outliers)
- Single bad measurement will throw off the mean
- Median = 50% quantile, does not care if you shift a point much further away —> robust against outliers
- Median is minimum of $X = \sum_i |\chi_i| = \sum_i |y_i - \text{median}|$
objective function

Estimates of spread

- Similar to mean, standard deviation is not robust against outliers
- Interquartile range [25%,75%] is robust, similar to the median (or any other range of quantiles near the center)
- Median absolute deviation:
$$\text{MAD} = \text{median} (|X_i - \text{median}(X_i)|)$$
- For Gaussian: std. dev. = 1.4826 MAD, generally useful to transform any robust range to its Gaussian equivalent