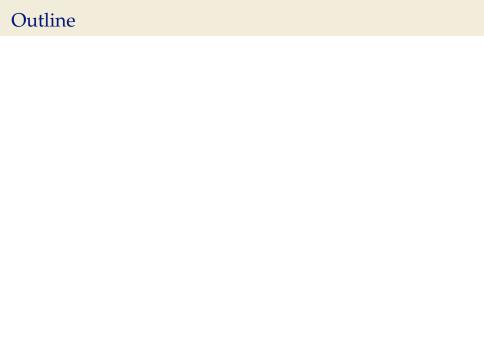
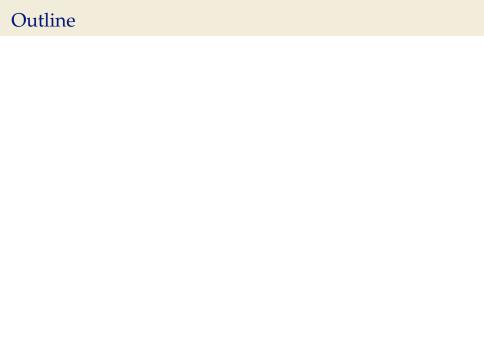
Regression and Probability

Neil D. Lawrence

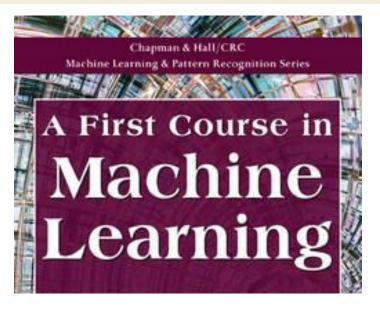
GPRS 11th February 2014



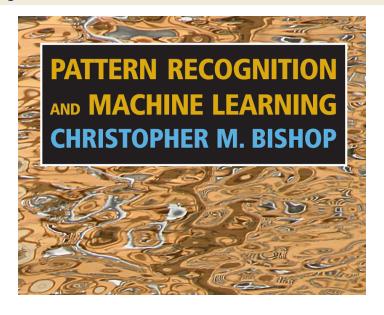


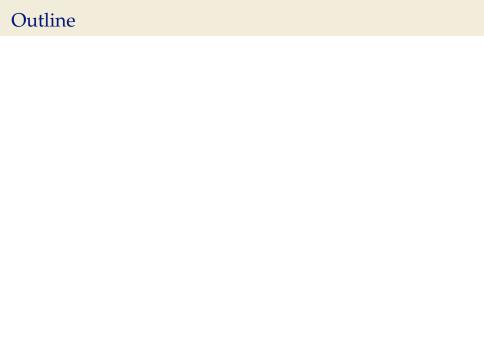


Rogers and Girolami



Bishop





Regression Examples

- ▶ Predict a real value, y_i given some inputs x_i .
- Predict quality of meat given spectral measurements (Tecator data).
- ▶ Radiocarbon dating, the C14 calibration curve: predict age given quantity of C14 isotope.
- Predict quality of different Go or Backgammon moves given expert rated training data.

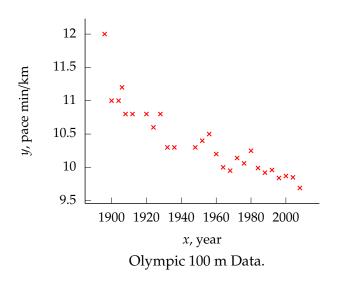
Olympic 100m Data

► Gold medal times for Olympic 100 m runners since 1896.



Image from Wikimedia Commons http://bit.ly/191adDC

Olympic 100m Data



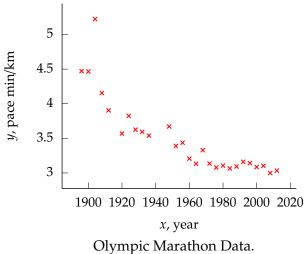
Olympic Marathon Data

- Gold medal times for Olympic Marathon since 1896.
- Marathons before 1924 didn't have a standardised distance.
- Present results using pace per km.
- In 1904 Marathon was badly organised leading to very slow times.



Image from Wikimedia Commons http://bit.ly/16kMKHQ

Olympic Marathon Data



data

 data: observations, could be actively or passively acquired (meta-data).

data +

 data: observations, could be actively or passively acquired (meta-data).

- data: observations, could be actively or passively acquired (meta-data).
- model: assumptions, based on previous experience (other data! transfer learning etc), or beliefs about the regularities of the universe. Inductive bias.

- data: observations, could be actively or passively acquired (meta-data).
- model: assumptions, based on previous experience (other data! transfer learning etc), or beliefs about the regularities of the universe. Inductive bias.

- data: observations, could be actively or passively acquired (meta-data).
- model: assumptions, based on previous experience (other data! transfer learning etc), or beliefs about the regularities of the universe. Inductive bias.
- prediction: an action to be taken or a categorization or a quality score.

$$y = mx + c$$

▶ **y**: winning time/pace.

$$y = mx + c$$

- ▶ **y**: winning time/pace.
- ► x: year of Olympics.

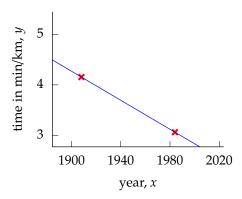
$$y = mx + c$$

- ▶ y: winning time/pace.
- ► x: year of Olympics.
- ▶ m: rate of improvement over time.

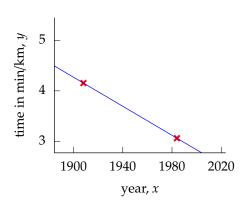
$$y = mx + c$$

- ▶ y: winning time/pace.
- ► x: year of Olympics.
- ▶ m: rate of improvement over time.
- c: winning time at year 0.

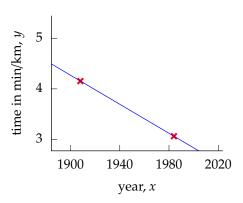
$$y_1 = mx_1 + c$$
$$y_2 = mx_2 + c$$



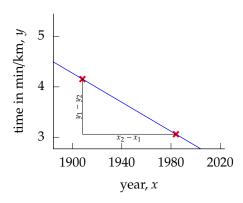
$$y_1 - y_2 = m(x_1 - x_2)$$



$$\frac{y_1 - y_2}{x_1 - x_2} = m$$



$$m = \frac{y_2 - y_1}{x_2 - x_1}$$
$$c = y_1 - mx_1$$

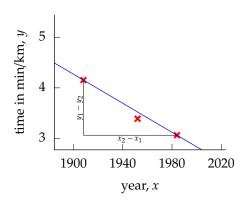


How do we deal with three simultaneous equations with only two unknowns?

$$y_1 = mx_1 + c$$

$$y_2 = mx_2 + c$$

$$y_3 = mx_3 + c$$



Overdetermined System

► With two unknowns and two observations:

$$y_1 = mx_1 + c$$
$$y_2 = mx_2 + c$$

Overdetermined System

▶ With two unknowns and two observations:

$$y_1 = mx_1 + c$$
$$y_2 = mx_2 + c$$

► Additional observation leads to *overdetermined* system.

$$y_3 = mx_3 + c$$

Overdetermined System

With two unknowns and two observations:

$$y_1 = mx_1 + c$$
$$y_2 = mx_2 + c$$

► Additional observation leads to *overdetermined* system.

$$y_3 = mx_3 + c$$

► This problem is solved through a noise model $\epsilon \sim \mathcal{N}(0, \sigma^2)$

$$y_1 = mx_1 + c + \epsilon_1$$

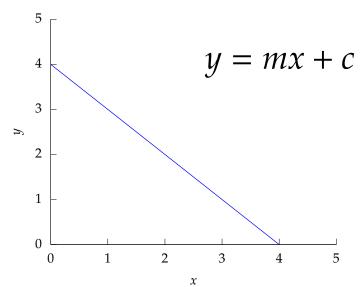
$$y_2 = mx_2 + c + \epsilon_2$$

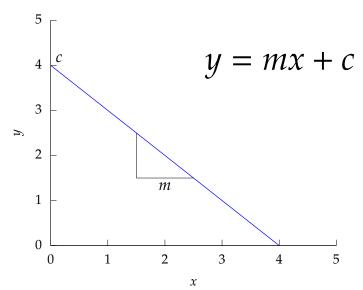
$$y_3 = mx_3 + c + \epsilon_3$$

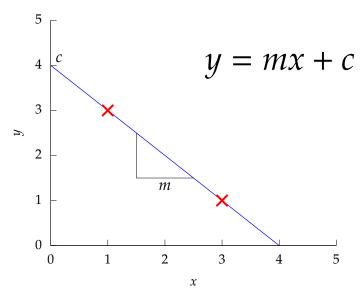
Noise Models

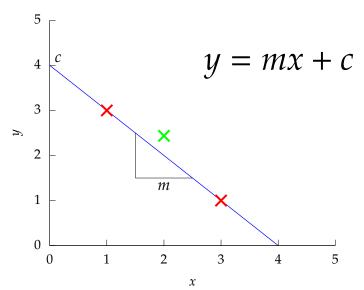
- We aren't modeling entire system.
- ▶ Noise model gives mismatch between model and data.
- Gaussian model justified by appeal to central limit theorem.
- ▶ Other models also possible (Student-*t* for heavy tails).
- Maximum likelihood with Gaussian noise leads to least squares.

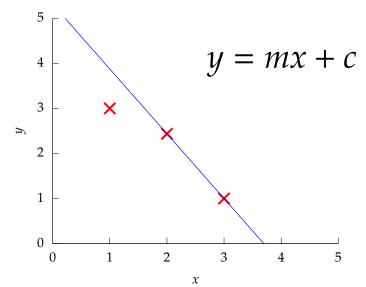
y = mx + c

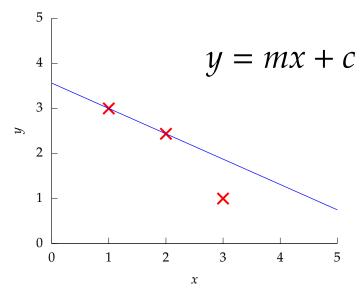


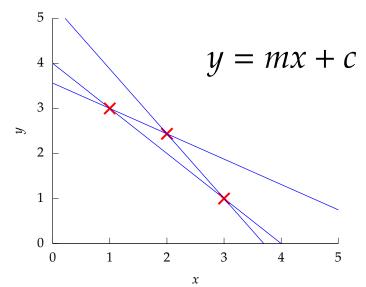












$$y = mx + c$$

point 1:
$$x = 1$$
, $y = 3$
 $3 = m + c$
point 2: $x = 3$, $y = 1$
 $1 = 3m + c$
point 3: $x = 2$, $y = 2.5$

2.5 = 2m + c

point 1:
$$x = 1$$
, $y = 3$
 $3 = m + c + \epsilon_1$
point 2: $x = 3$, $y = 1$
 $1 = 3m + c + \epsilon_2$
point 3: $x = 2$, $y = 2.5$

$$2.5 = 2m + c + \epsilon_3$$

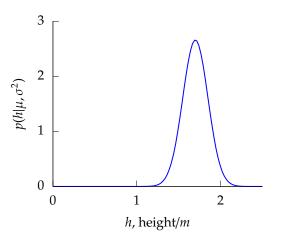
The Gaussian Density

▶ Perhaps the most common probability density.

$$p(y|\mu,\sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y-\mu)^2}{2\sigma^2}\right)$$
$$\stackrel{\triangle}{=} \mathcal{N}\left(y|\mu,\sigma^2\right)$$

► The Gaussian density.

Gaussian Density



The Gaussian PDF with $\mu = 1.7$ and variance $\sigma^2 = 0.0225$. Mean shown as red line. It could represent the heights of a population of students.

Gaussian Density

$$\mathcal{N}(y|\mu,\sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y-\mu)^2}{2\sigma^2}\right)$$

 σ^2 is the variance of the density and μ is the mean.

Sum of Gaussians

▶ Sum of Gaussian variables is also Gaussian.

$$y_i \sim \mathcal{N}\left(\mu_i, \sigma_i^2\right)$$

Sum of Gaussians

▶ Sum of Gaussian variables is also Gaussian.

$$y_i \sim \mathcal{N}\left(\mu_i, \sigma_i^2\right)$$

And the sum is distributed as

$$\sum_{i=1}^{n} y_i \sim \mathcal{N}\left(\sum_{i=1}^{n} \mu_i, \sum_{i=1}^{n} \sigma_i^2\right)$$

Sum of Gaussians

▶ Sum of Gaussian variables is also Gaussian.

$$y_i \sim \mathcal{N}\left(\mu_i, \sigma_i^2\right)$$

And the sum is distributed as

$$\sum_{i=1}^{n} y_i \sim \mathcal{N}\left(\sum_{i=1}^{n} \mu_i, \sum_{i=1}^{n} \sigma_i^2\right)$$

(*Aside*: As sum increases, sum of non-Gaussian, finite variance variables is also Gaussian [central limit theorem].)

Sum of Gaussians

▶ Sum of Gaussian variables is also Gaussian.

$$y_i \sim \mathcal{N}\left(\mu_i, \sigma_i^2\right)$$

And the sum is distributed as

$$\sum_{i=1}^{n} y_i \sim \mathcal{N}\left(\sum_{i=1}^{n} \mu_i, \sum_{i=1}^{n} \sigma_i^2\right)$$

(*Aside*: As sum increases, sum of non-Gaussian, finite variance variables is also Gaussian [central limit theorem].)

Scaling a Gaussian

► Scaling a Gaussian leads to a Gaussian.

Scaling a Gaussian

► Scaling a Gaussian leads to a Gaussian.

$$y \sim \mathcal{N}\left(\mu, \sigma^2\right)$$

Scaling a Gaussian

► Scaling a Gaussian leads to a Gaussian.

$$y \sim \mathcal{N}\left(\mu, \sigma^2\right)$$

And the scaled density is distributed as

$$wy \sim \mathcal{N}\left(w\mu, w^2\sigma^2\right)$$

A Probabilistic Process

▶ Set the mean of Gaussian to be a function.

$$p(y_i|x_i) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y_i - f(x_i))^2}{2\sigma^2}\right).$$

- ► This gives us a 'noisy function'.
- ► This is known as a process.

Height as a Function of Weight

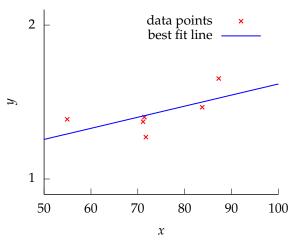
- ► In the standard Gaussian, parametized by mean and variance.
- ▶ Make the mean a linear function of an *input*.
- ► This leads to a regression model.

$$y_i = f(x_i) + \epsilon_i,$$

 $\epsilon_i \sim \mathcal{N}(0, \sigma^2).$

▶ Assume y_i is height and x_i is weight.

Linear Function



A linear regression between x and y.

Data Point Likelihood

Likelihood of an individual data point

$$p(y_i|x_i, m, c) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y_i - mx_i - c)^2}{2\sigma^2}\right).$$

▶ Parameters are gradient, m, offset, c of the function and noise variance σ^2 .

- ▶ If the noise, ϵ_i is sampled independently for each data point.
- ► Each data point is independent (given *m* and *c*).
- ► For independent variables:

$$p(\mathbf{y}) = \prod_{i=1}^n p(y_i)$$

- ▶ If the noise, ϵ_i is sampled independently for each data point.
- ► Each data point is independent (given *m* and *c*).
- ► For independent variables:

$$p(\mathbf{y}|\mathbf{x},m,c) = \prod_{i=1}^{n} p(y_i|x_i,m,c)$$

- ▶ If the noise, ϵ_i is sampled independently for each data point.
- ► Each data point is independent (given *m* and *c*).
- ► For independent variables:

$$p(\mathbf{y}|\mathbf{x}, m, c) = \prod_{i=1}^{n} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y_i - mx_i - c)^2}{2\sigma^2}\right).$$

- ▶ If the noise, ϵ_i is sampled independently for each data point.
- ► Each data point is independent (given *m* and *c*).
- ► For independent variables:

$$p(\mathbf{y}|\mathbf{x}, m, c) = \frac{1}{(2\pi\sigma^2)^{\frac{n}{2}}} \exp\left(-\frac{\sum_{i=1}^{n} (y_i - mx_i - c)^2}{2\sigma^2}\right).$$

Log Likelihood Function

► Normally work with the log likelihood:

$$L(m,c,\sigma^2) = -\frac{n}{2}\log 2\pi - \frac{n}{2}\log \sigma^2 - \sum_{i=1}^n \frac{(y_i - mx_i - c)^2}{2\sigma^2}.$$

Consistency of Maximum Likelihood

- If data was really generated according to probability we specified.
- ► Correct parameters will be recovered in limit as $n \to \infty$.
- ► This can be proven through sample based approximations (law of large numbers) of "KL divergences".
- Mainstay of classical statistics.

Probabilistic Interpretation of the Error Function

- Probabilistic Interpretation for Error Function is Negative Log Likelihood.
- ► *Minimizing* error function is equivalent to *maximizing* log likelihood.
- Maximizing log likelihood is equivalent to maximizing the likelihood because log is monotonic.
- Probabilistic interpretation: Minimizing error function is equivalent to maximum likelihood with respect to parameters.

Error Function

 Negative log likelihood is the error function leading to an error function

$$E(m, c, \sigma^2) = \frac{n}{2} \log \sigma^2 + \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - mx_i - c)^2.$$

 Learning proceeds by minimizing this error function for the data set provided.

Connection: Sum of Squares Error

▶ Ignoring terms which don't depend on *m* and *c* gives

$$E(m,c) \propto \sum_{i=1}^{n} (y_i - f(x_i))^2$$

where $f(x_i) = mx_i + c$.

- ► This is known as the *sum of squares* error function.
- Commonly used and is closely associated with the Gaussian likelihood.

Mathematical Interpretation

- What is the mathematical interpretation?
 - ► There is a cost function.
 - ► It expresses mismatch between your prediction and reality.

$$E(\mathbf{w}) = \sum_{i=1}^{n} (y_i - mx_i - c)^2$$

► This is known as the sum of squares error.

- ► Learning is minimization of the cost function.
- ▶ At the minima the gradient is zero.
- Coordinate ascent, find gradient in each coordinate and set to zero.

$$\frac{\mathrm{d}E(m)}{\mathrm{d}m} = -2\sum_{i=1}^{n} x_i \left(y_i - mx_i - c\right)$$

- ► Learning is minimization of the cost function.
- ▶ At the minima the gradient is zero.
- Coordinate ascent, find gradient in each coordinate and set to zero.

$$0 = -2\sum_{i=1}^{n} x_i (y_i - mx_i - c)$$

- ► Learning is minimization of the cost function.
- ▶ At the minima the gradient is zero.
- Coordinate ascent, find gradient in each coordinate and set to zero.

$$0 = -2\sum_{i=1}^{n} x_i y_i + 2\sum_{i=1}^{n} mx_i^2 + 2\sum_{i=1}^{n} cx_i$$

- ▶ Learning is minimization of the cost function.
- ▶ At the minima the gradient is zero.
- Coordinate ascent, find gradient in each coordinate and set to zero.

$$m = \frac{\sum_{i=1}^{n} (y_i - c) x_i}{\sum_{i=1}^{n} x_i^2}$$

- ► Learning is minimization of the cost function.
- ▶ At the minima the gradient is zero.
- Coordinate ascent, find gradient in each coordinate and set to zero.

$$\frac{\mathrm{d}E(c)}{\mathrm{d}c} = -2\sum_{i=1}^{n} (y_i - mx_i - c)$$

- ► Learning is minimization of the cost function.
- ▶ At the minima the gradient is zero.
- Coordinate ascent, find gradient in each coordinate and set to zero.

$$0 = -2\sum_{i=1}^{n} (y_i - mx_i - c)$$

- ► Learning is minimization of the cost function.
- ▶ At the minima the gradient is zero.
- Coordinate ascent, find gradient in each coordinate and set to zero.

$$0 = -2\sum_{i=1}^{n} y_i + 2\sum_{i=1}^{n} mx_i + 2nc$$

- ► Learning is minimization of the cost function.
- ▶ At the minima the gradient is zero.
- Coordinate ascent, find gradient in each coordinate and set to zero.

$$c = \frac{\sum_{i=1}^{n} (y_i - mx_i)}{n}$$

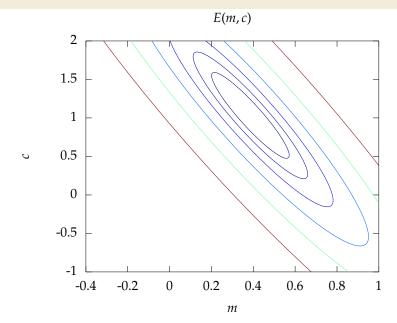
Fixed Point Updates

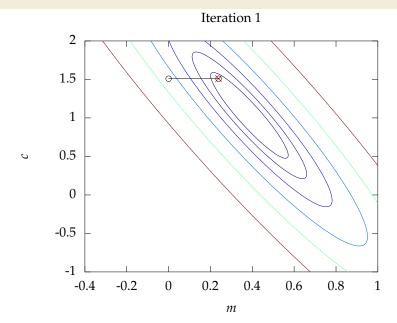
Worked example.

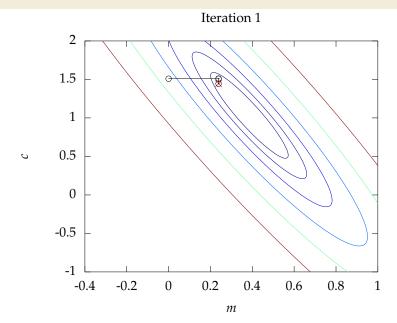
$$c^* = \frac{\sum_{i=1}^{n} (y_i - m^* x_i)}{n},$$

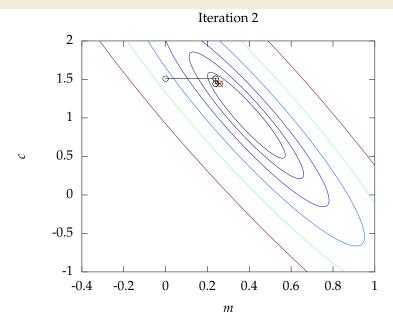
$$m^* = \frac{\sum_{i=1}^{n} x_i (y_i - c^*)}{\sum_{i=1}^{n} x_i^2},$$

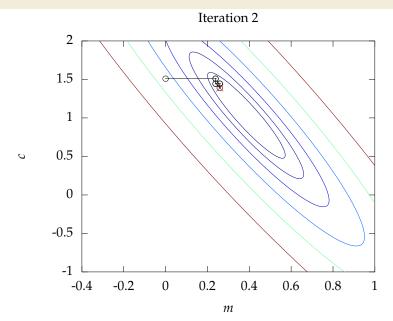
$$\sigma^{2^*} = \frac{\sum_{i=1}^{n} (y_i - m^* x_i - c^*)^2}{n}$$

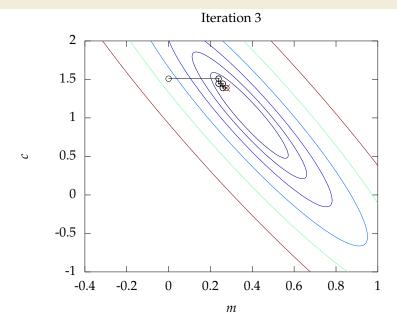


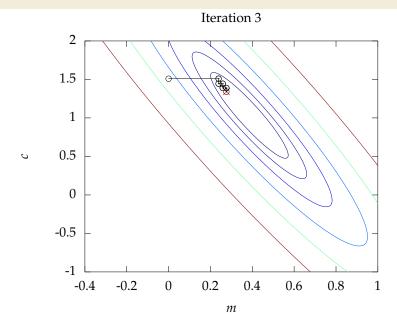


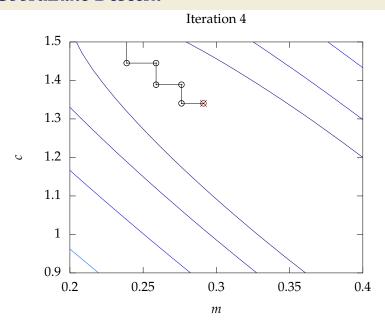


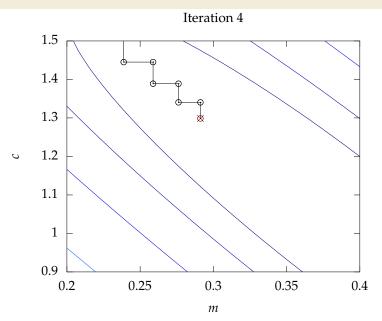


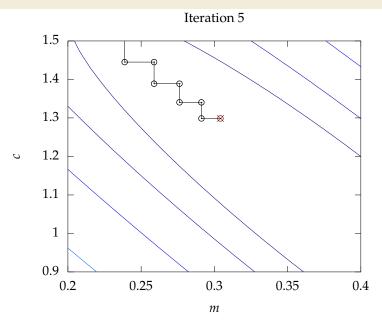


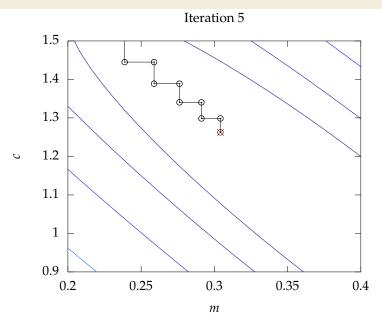


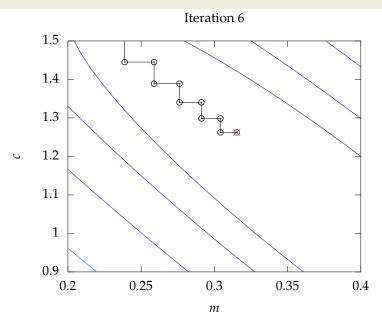


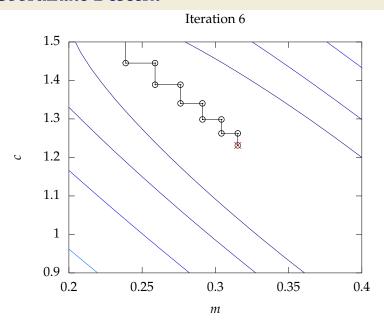


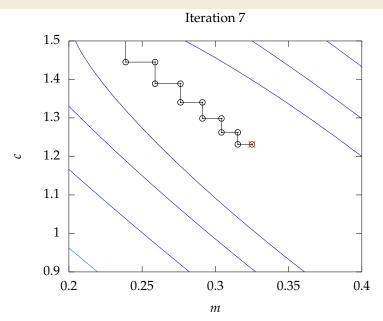


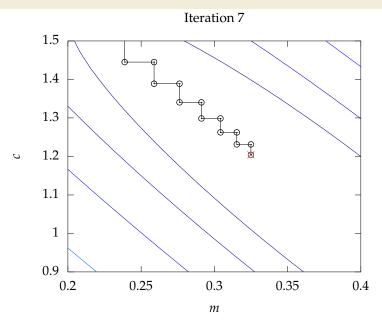


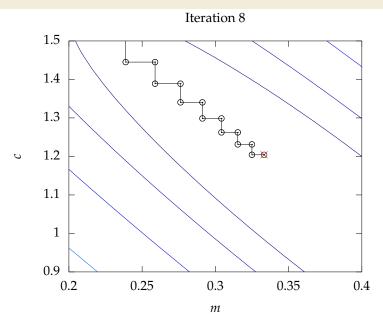


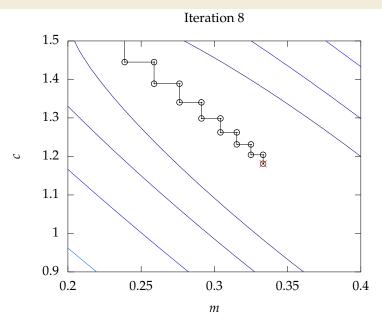


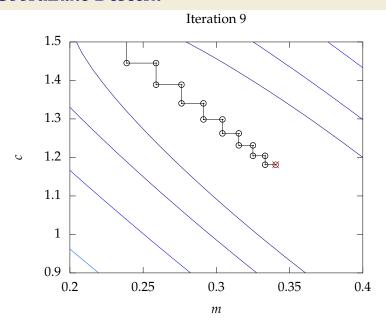


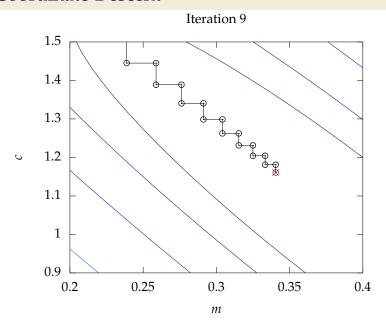


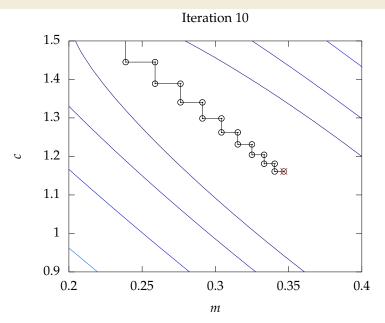


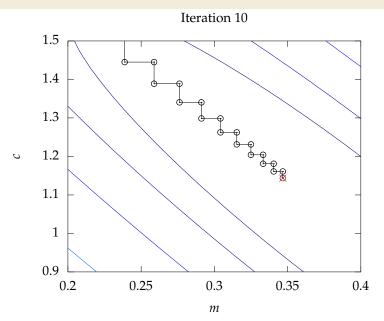


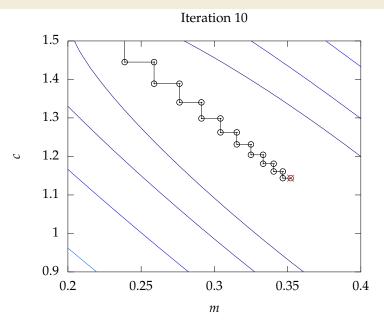


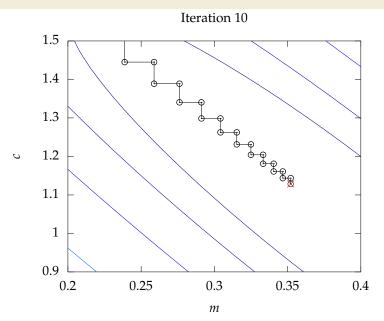


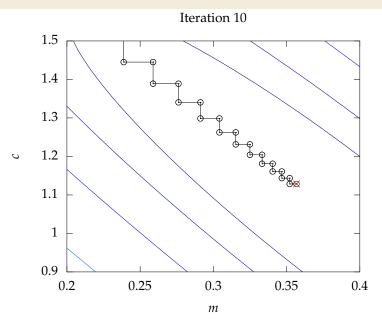


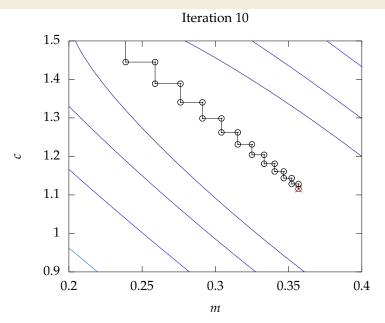


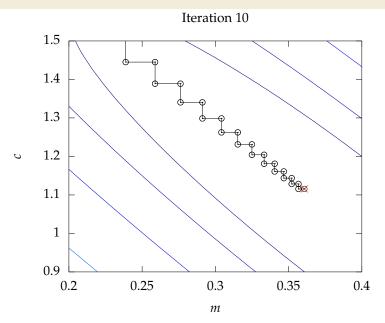


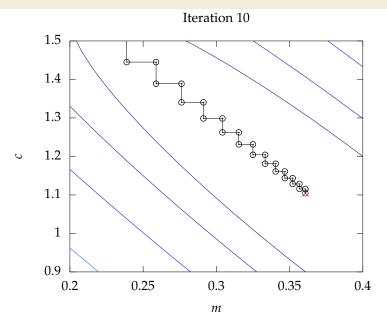


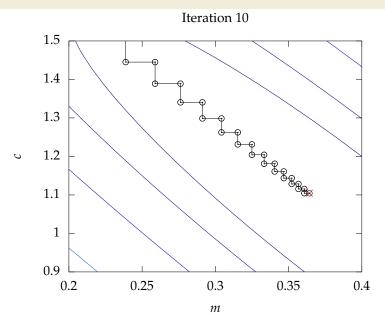


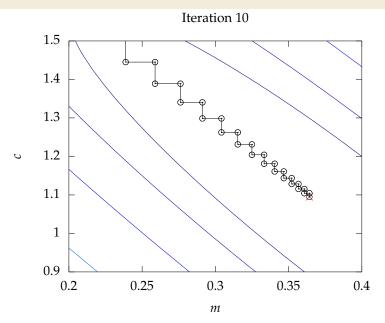


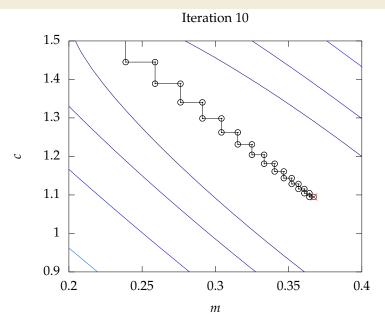


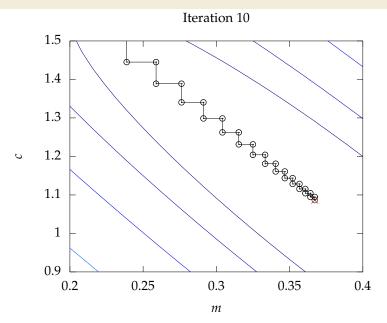


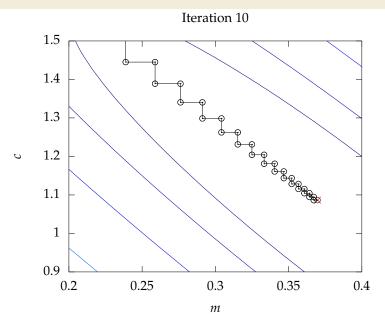


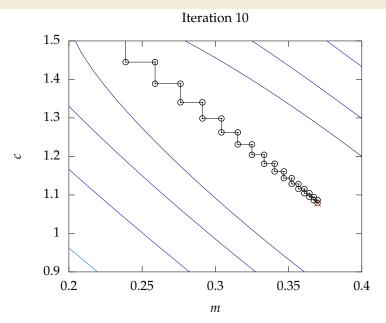


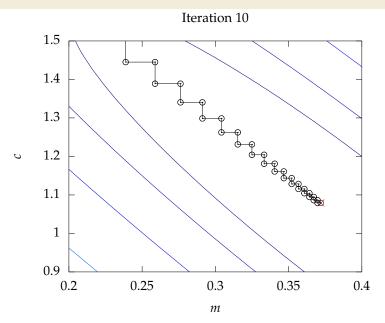


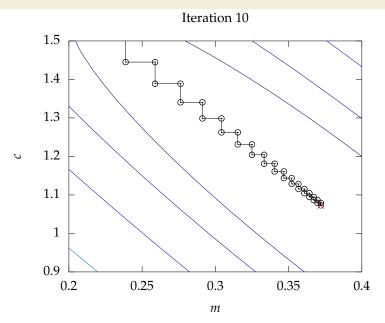


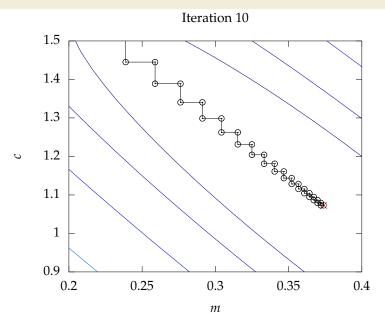


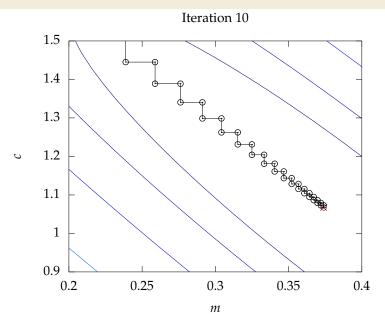


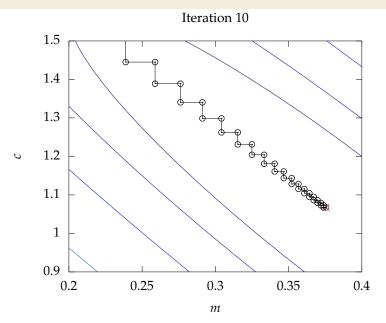


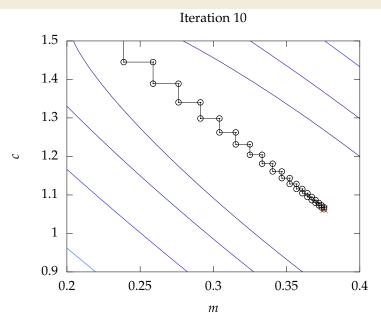


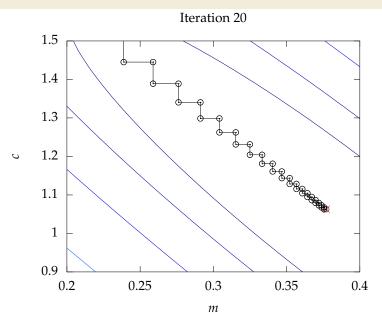


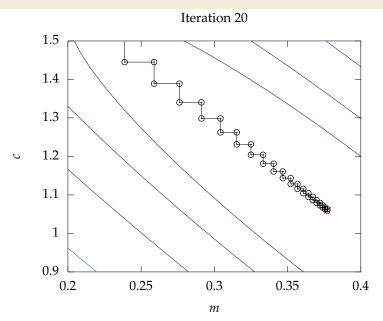


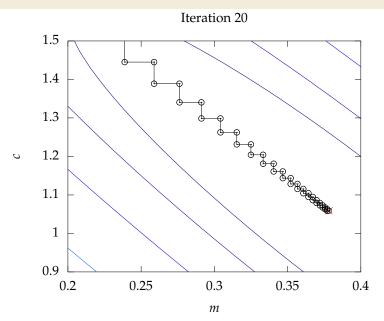


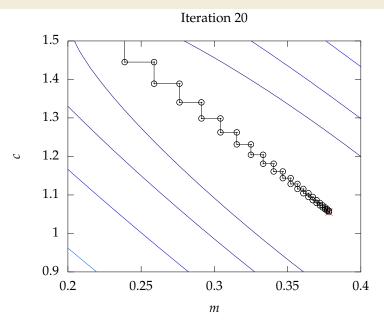


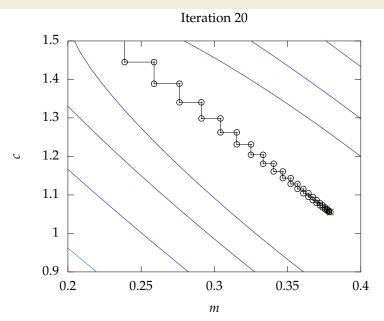


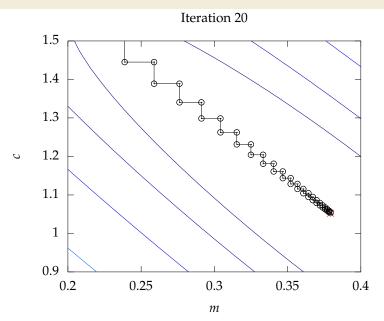


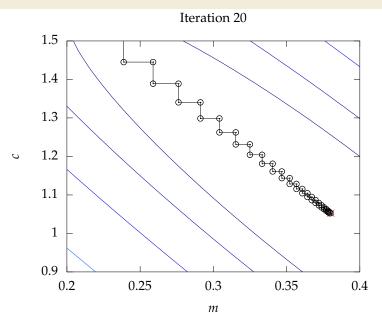


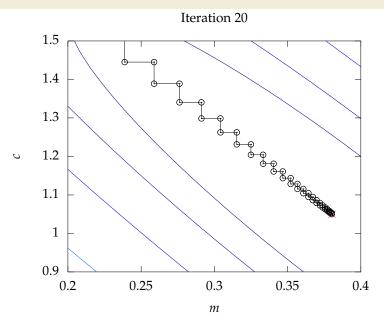


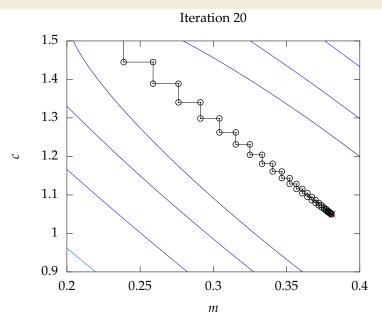


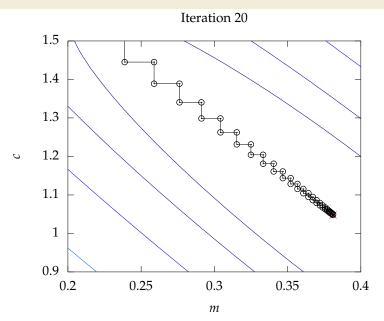


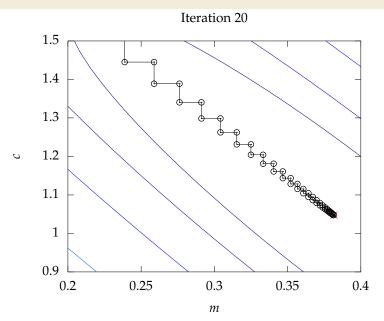


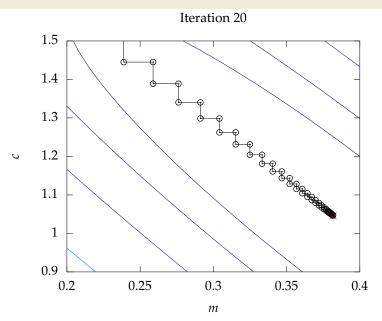


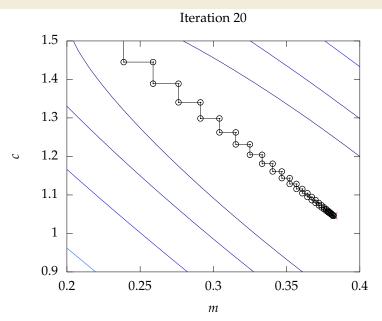


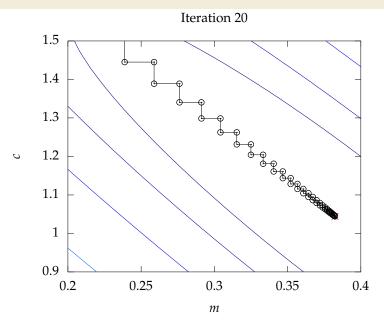


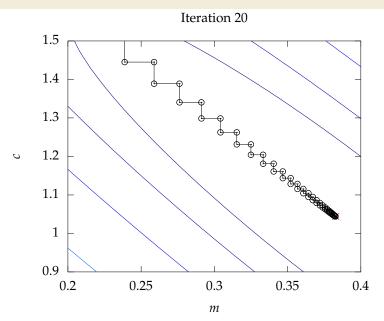


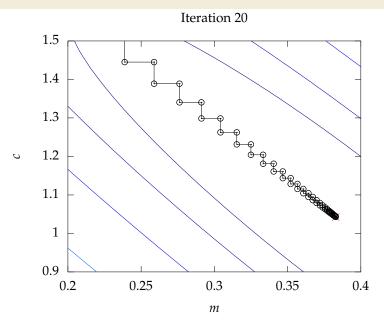


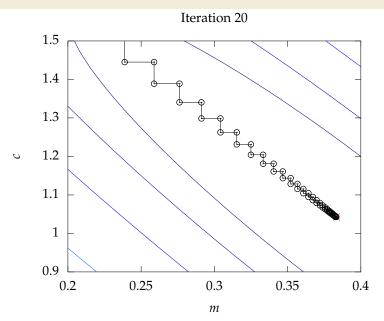


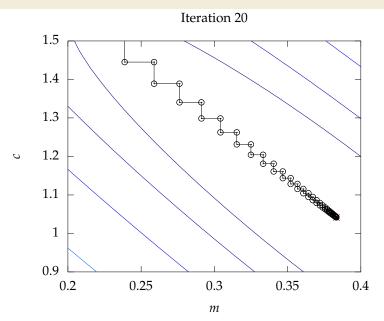


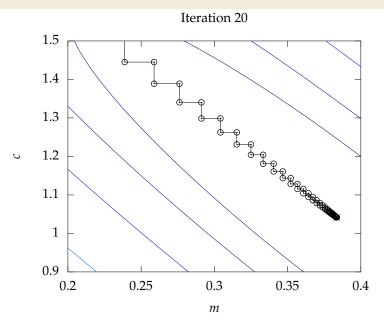


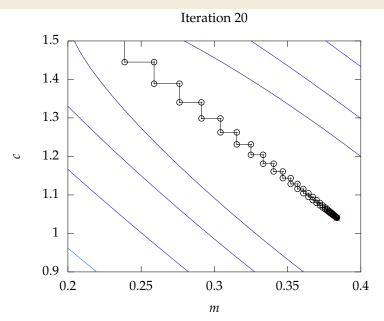


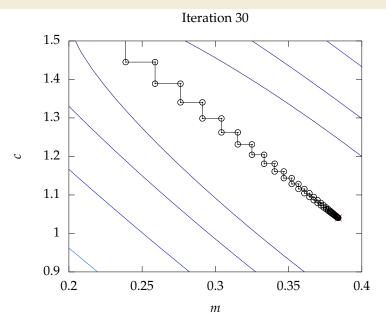


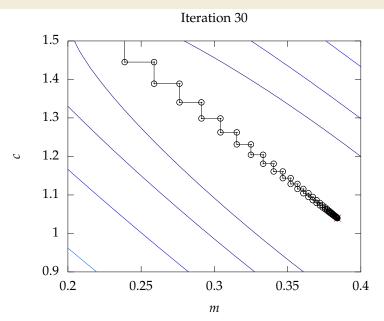


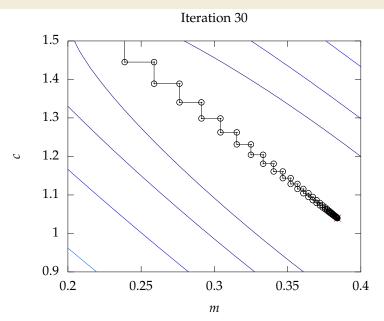








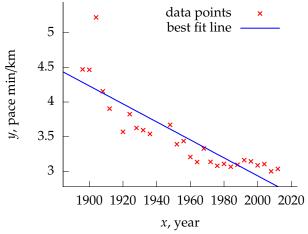




Important Concepts Not Covered

- Optimization methods.
 - Second order methods, conjugate gradient, quasi-Newton and Newton.
 - Effective heuristics such as momentum.
- ► Local vs global solutions.

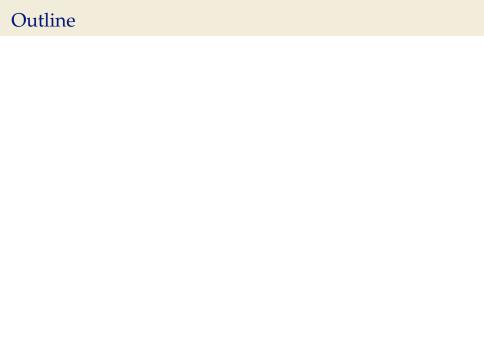
Linear Function



Linear regression for Male Olympics Marathon Gold Medal times.

Reading

- ► Section 1.2.5 of Bishop up to equation 1.65.
- ► Section 1.1-1.2 of Rogers and Girolami for fitting linear models.



Basis Functions

Nonlinear Regression

- ► Problem with Linear Regression—x may not be linearly related to y.
- ▶ Potential solution: create a feature space: define $\phi(\mathbf{x})$ where $\phi(\cdot)$ is a nonlinear function of \mathbf{x} .
- Model for target is a linear combination of these nonlinear functions

$$f(\mathbf{x}) = \sum_{i=1}^{K} w_i \phi_i(\mathbf{x})$$
 (1)

Quadratic Basis

▶ Basis functions can be global. E.g. quadratic basis:

$$[1,x,x^2]$$

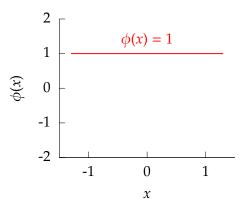


Figure : A quadratic basis.

Quadratic Basis

▶ Basis functions can be global. E.g. quadratic basis:

$$[1, x, x^2]$$

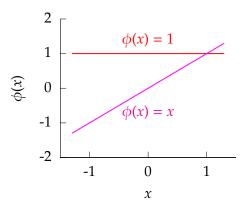


Figure : A quadratic basis.

Quadratic Basis

▶ Basis functions can be global. E.g. quadratic basis:

$$[1, x, x^2]$$

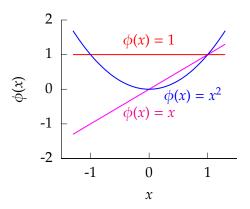


Figure : A quadratic basis.

Functions Derived from Quadratic Basis

$$f(x) = w_1 + w_2 x + w_3 x^2$$

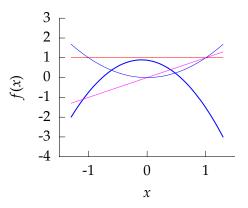


Figure : Function from quadratic basis with weights $w_1 = 0.87466$, $w_2 = -0.38835$, $w_3 = -2.0058$.

Functions Derived from Quadratic Basis

$$f(x) = w_1 + w_2 x + w_3 x^2$$

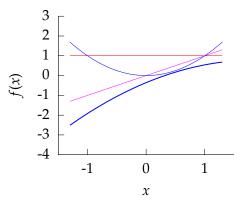


Figure : Function from quadratic basis with weights $w_1 = -0.35908$, $w_2 = 1.2274$, $w_3 = -0.32825$.

Functions Derived from Quadratic Basis

$$f(x) = w_1 + w_2 x + w_3 x^2$$

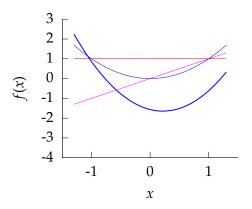


Figure : Function from quadratic basis with weights $w_1 = -1.5638$, $w_2 = -0.73577$, $w_3 = 1.6861$.

Radial Basis Functions

► Or they can be local. E.g. radial (or Gaussian) basis $\phi_j(x) = \exp\left(-\frac{(x-\mu_j)^2}{\ell^2}\right)$

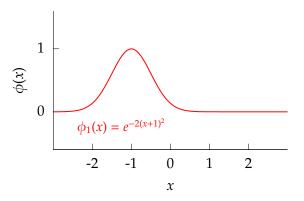


Figure: Radial basis functions.

Radial Basis Functions

► Or they can be local. E.g. radial (or Gaussian) basis $\phi_j(x) = \exp\left(-\frac{(x-\mu_j)^2}{\ell^2}\right)$

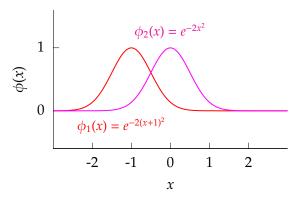


Figure: Radial basis functions.

Radial Basis Functions

► Or they can be local. E.g. radial (or Gaussian) basis $\phi_j(x) = \exp\left(-\frac{(x-\mu_j)^2}{\ell^2}\right)$

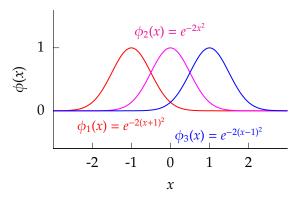


Figure: Radial basis functions.

Functions Derived from Radial Basis

$$f(x) = w_1 e^{-2(x+1)^2} + w_2 e^{-2x^2} + w_3 e^{-2(x-1)^2}$$

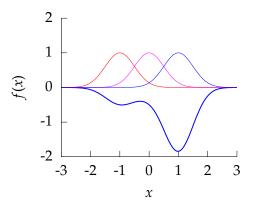


Figure : Function from radial basis with weights $w_1 = -0.47518$, $w_2 = -0.18924$, $w_3 = -1.8183$.

Functions Derived from Radial Basis

$$f(x) = w_1 e^{-2(x+1)^2} + w_2 e^{-2x^2} + w_3 e^{-2(x-1)^2}$$

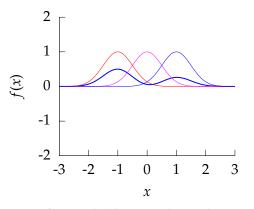


Figure : Function from radial basis with weights $w_1 = 0.50596$, $w_2 = -0.046315$, $w_3 = 0.26813$.

Functions Derived from Radial Basis

$$f(x) = w_1 e^{-2(x+1)^2} + w_2 e^{-2x^2} + w_3 e^{-2(x-1)^2}$$

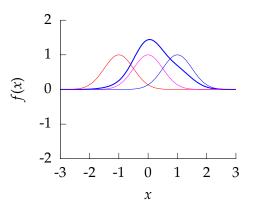


Figure : Function from radial basis with weights $w_1 = 0.07179$, $w_2 = 1.3591$, $w_3 = 0.50604$.

Reading

- ► Chapter 1, pg 1-6 of Bishop.
- ► Section 1.4 of Rogers and Girolami.
- ► Chapter 3, Section 3.1 of Bishop up to pg 143.

Multi-dimensional Inputs

- ► Multivariate functions involve more than one input.
- ► Height might be a function of weight and gender.
- ▶ There could be other contributory factors.
- ▶ Place these factors in a feature vector \mathbf{x}_i .
- Linear function is now defined as

$$f(\mathbf{x}_i) = \sum_{j=1}^{q} w_j x_{i,j} + c$$

Vector Notation

mo

► Write in vector notation,

$$f(\mathbf{x}_i) = \mathbf{w}^{\top} \mathbf{x}_i + c$$

► Can absorb c into \mathbf{w} by assuming extra input x_0 which is always 1.

$$f(\mathbf{x}_i) = \mathbf{w}^{\top} \mathbf{x}_i$$

Log Likelihood for Multivariate Regression

► The likelihood of a single data point is

$$p\left(y_i|x_i\right) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{\left(y_i - \mathbf{w}^\top \mathbf{x}_i\right)^2}{2\sigma^2}\right).$$

Leading to a log likelihood for the data set of

$$L(\mathbf{w}, \sigma^2) = -\frac{n}{2} \log \sigma^2 - \frac{n}{2} \log 2\pi - \frac{\sum_{i=1}^n (y_i - \mathbf{w}^\top \mathbf{x}_i)^2}{2\sigma^2}.$$

And a corresponding error function of

$$E(\mathbf{w}, \sigma^2) = \frac{n}{2} \log \sigma^2 + \frac{\sum_{i=1}^n (y_i - \mathbf{w}^\top \mathbf{x}_i)^2}{2\sigma^2}.$$

Expand the Brackets

$$E(\mathbf{w}, \sigma^2) = \frac{n}{2} \log \sigma^2 + \frac{1}{2\sigma^2} \sum_{i=1}^n y_i^2 - \frac{1}{\sigma^2} \sum_{i=1}^n y_i \mathbf{w}^\top \mathbf{x}_i$$

$$+ \frac{1}{2\sigma^2} \sum_{i=1}^n \mathbf{w}^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{w} + \text{const.}$$

$$= \frac{n}{2} \log \sigma^2 + \frac{1}{2\sigma^2} \sum_{i=1}^n y_i^2 - \frac{1}{\sigma^2} \mathbf{w}^\top \sum_{i=1}^n \mathbf{x}_i y_i$$

$$+ \frac{1}{2\sigma^2} \mathbf{w}^\top \left[\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top \right] \mathbf{w} + \text{const.}$$

Multivariate Derivatives

- ▶ We will need some multivariate calculus.
- ► For now some simple multivariate differentiation:

$$\frac{\mathrm{d}\mathbf{a}^{\top}\mathbf{w}}{\mathrm{d}\mathbf{w}} = \mathbf{a}$$

and

$$\frac{\mathrm{d}\mathbf{w}^{\top}\mathbf{A}\mathbf{w}}{\mathrm{d}\mathbf{w}} = \left(\mathbf{A} + \mathbf{A}^{\top}\right)\mathbf{w}$$

or if **A** is symmetric (*i.e.* $\mathbf{A} = \mathbf{A}^{\top}$)

$$\frac{\mathrm{d}\mathbf{w}^{\top}\mathbf{A}\mathbf{w}}{\mathrm{d}\mathbf{w}} = 2\mathbf{A}\mathbf{w}.$$

Differentiate

Differentiating with respect to the vector \mathbf{w} we obtain

$$\frac{\partial L(\mathbf{w}, \beta)}{\partial \mathbf{w}} = \beta \sum_{i=1}^{n} \mathbf{x}_{i} y_{i} - \beta \left[\sum_{i=1}^{n} \mathbf{x}_{i} \mathbf{x}_{i}^{\mathsf{T}} \right] \mathbf{w}$$

Leading to

$$\mathbf{w}^* = \left[\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top\right]^{-1} \sum_{i=1}^n \mathbf{x}_i y_i,$$

Rewrite in matrix notation:

$$\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top = \mathbf{X}^\top \mathbf{X}$$

$$\sum_{i=1}^{n} \mathbf{x}_i y_i = \mathbf{X}^{\top} \mathbf{y}$$

Update Equations

▶ Update for w*.

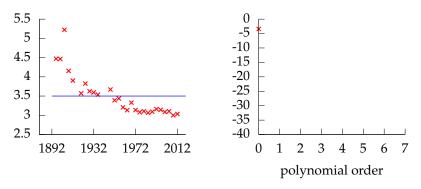
$$\mathbf{w}^* = \left(\mathbf{X}^\top \mathbf{X}\right)^{-1} \mathbf{X}^\top \mathbf{y}$$

► The equation for σ^{2^*} may also be found

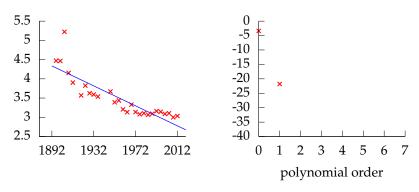
$$\sigma^{2^*} = \frac{\sum_{i=1}^n \left(y_i - \mathbf{w}^{* \top} \mathbf{x}_i \right)^2}{n}.$$

Reading

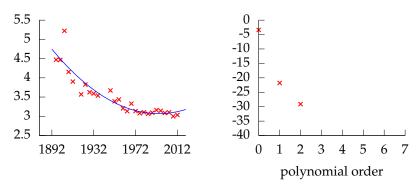
Section 1.3 of Rogers and Girolami for Matrix & Vector Review.



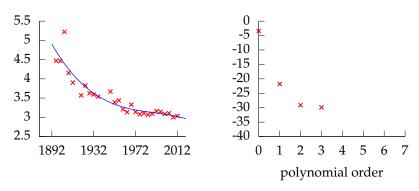
Left: fit to data, *Right*: model error. Polynomial order 0, model error -3.3989, $\sigma^2 = 0.286$, $\sigma = 0.535$.



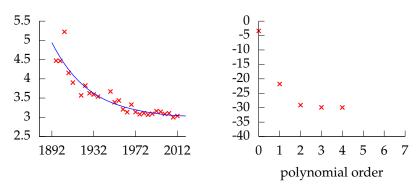
Left: fit to data, *Right*: model error. Polynomial order 1, model error -21.772, $\sigma^2 = 0.0733$, $\sigma = 0.271$.



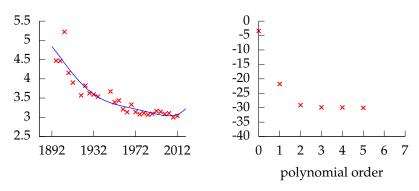
Left: fit to data, *Right*: model error. Polynomial order 2, model error -29.101, $\sigma^2 = 0.0426$, $\sigma = 0.206$.



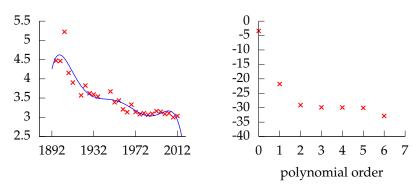
Left: fit to data, *Right*: model error. Polynomial order 3, model error -29.907, $\sigma^2 = 0.0401$, $\sigma = 0.200$.



Left: fit to data, *Right*: model error. Polynomial order 4, model error -29.943, $\sigma^2 = 0.0400$, $\sigma = 0.200$.



Left: fit to data, *Right*: model error. Polynomial order 5, model error -30.056, $\sigma^2 = 0.0397$, $\sigma = 0.199$.



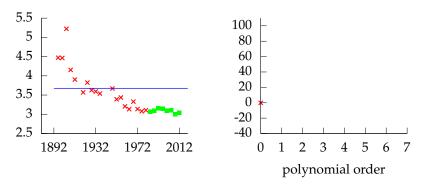
Left: fit to data, *Right*: model error. Polynomial order 6, model error -32.866, $\sigma^2 = 0.0322$, $\sigma = 0.180$.

Overfitting

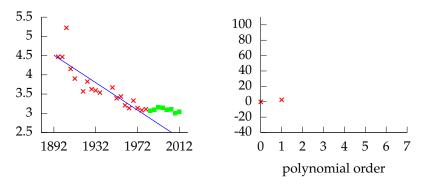
- ► Increase number of basis functions, we obtain a better 'fit' to the data.
- ► How will the model perform on previously unseen data?

Training and Test Sets

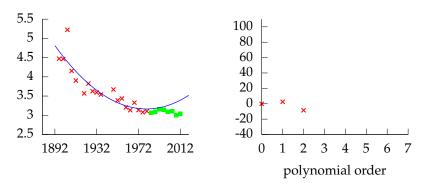
- We call the data used for fitting the model the 'training set'.
- ▶ Data not used for training, but when the model is applied 'in the field' is called the 'test data'.
- Challenge for generalization is to ensure a good performance on test data given only training data.



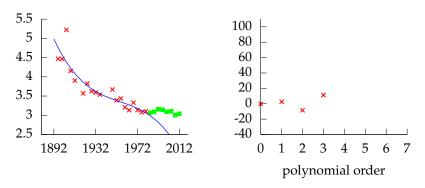
Left: fit to data, *Right*: model error. Polynomial order 0, training error -1.8774, validation error -0.13132, $\sigma^2 = 0.302$, $\sigma = 0.549$.



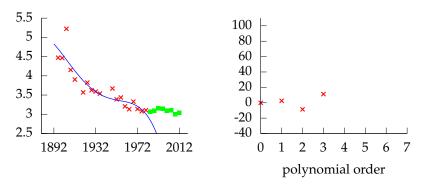
Left: fit to data, *Right*: model error. Polynomial order 1, training error -15.325, validation error 2.5863, $\sigma^2 = 0.0733$, $\sigma = 0.271$.



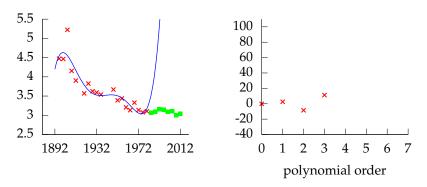
Left: fit to data, *Right*: model error. Polynomial order 2, training error -17.579, validation error -8.4831, $\sigma^2 = 0.0578$, $\sigma = 0.240$.



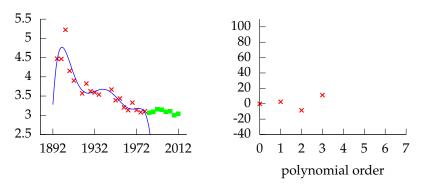
Left: fit to data, *Right*: model error. Polynomial order 3, training error -18.064, validation error 11.27, $\sigma^2 = 0.0549$, $\sigma = 0.234$.



Left: fit to data, *Right*: model error. Polynomial order 4, training error -18.245, validation error 232.92, $\sigma^2 = 0.0539$, $\sigma = 0.232$.

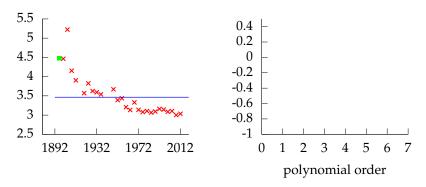


Left: fit to data, *Right*: model error. Polynomial order 5, training error -20.471, validation error 9898.1, $\sigma^2 = 0.0426$, $\sigma = 0.207$.

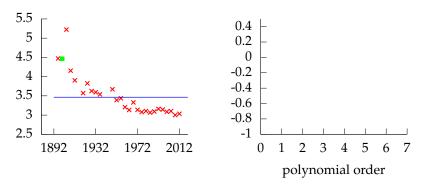


Left: fit to data, *Right*: model error. Polynomial order 6, training error -22.881, validation error 67775, $\sigma^2 = 0.0331$, $\sigma = 0.182$.

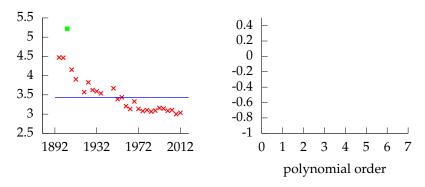
- ► Take training set and remove one point.
- ► Train on the remaining data.
- Compute the error on the point you removed (which wasn't in the training data).
- ▶ Do this for each point in the training set in turn.
- Average the resulting error. This is the leave one out error.



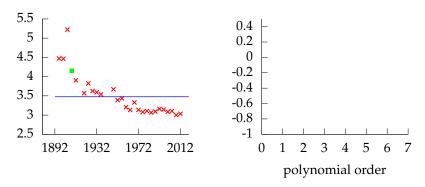
Polynomial order 0, training error -3.346, leave one out error 0.045811.

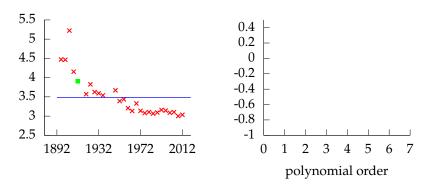


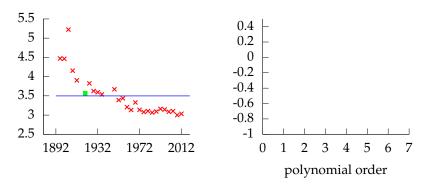
Polynomial order 0, training error -3.346, leave one out error 0.045811.

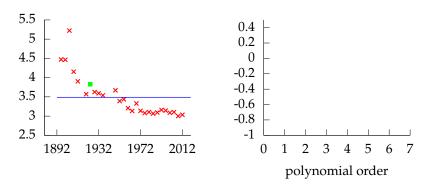


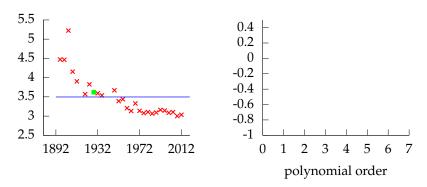
Polynomial order 0, training error -3.346, leave one out error 0.045811.



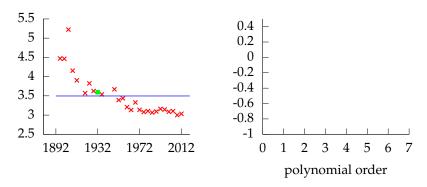


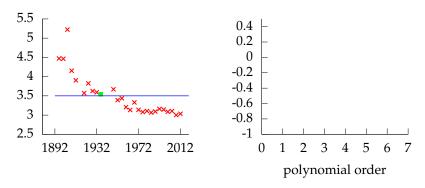


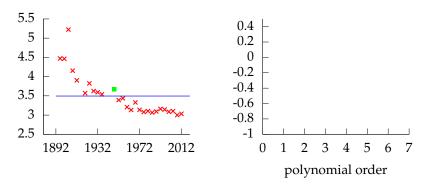


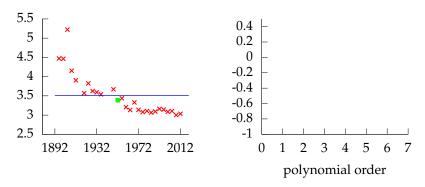


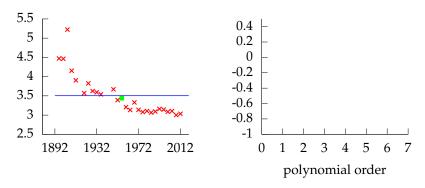
Polynomial order 0, training error -3.346, leave one out error 0.045811.

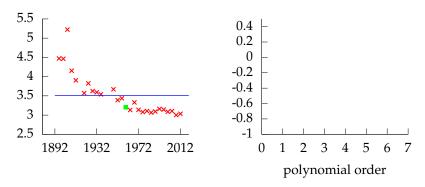


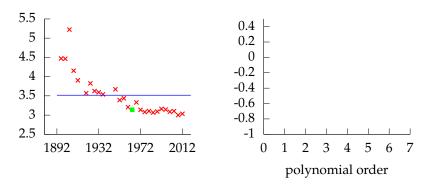


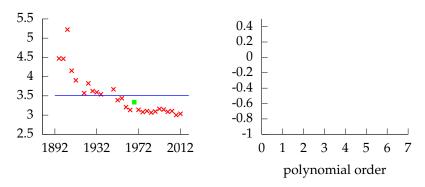


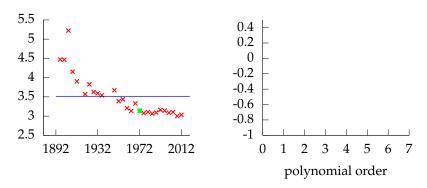


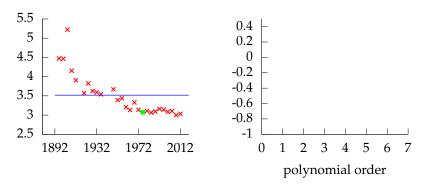


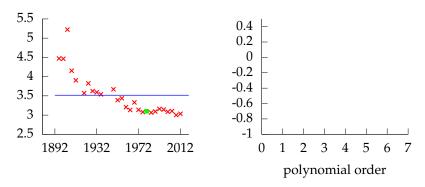


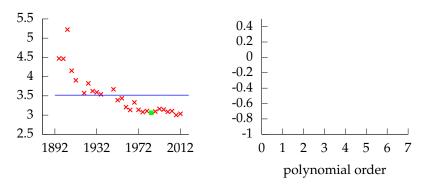


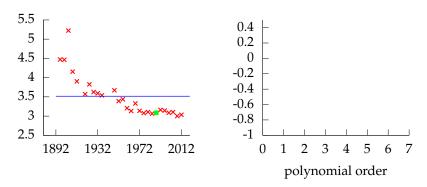




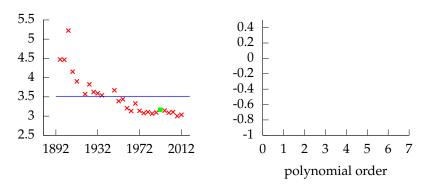




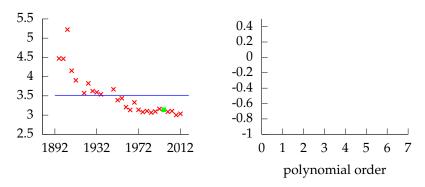


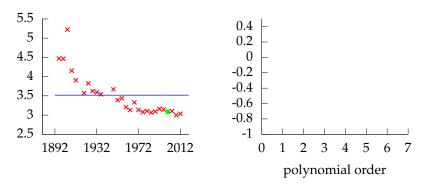


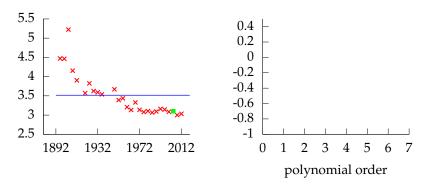
Polynomial order 0, training error -3.346, leave one out error 0.045811.

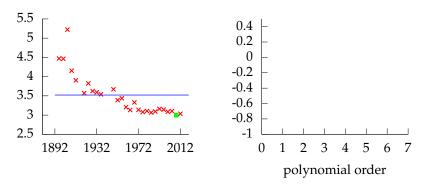


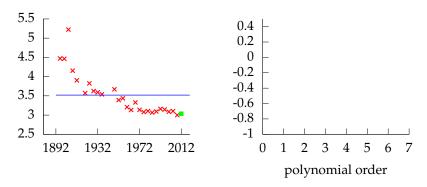
Polynomial order 0, training error -3.346, leave one out error 0.045811.

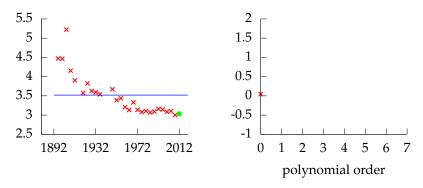


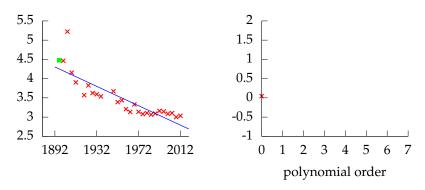




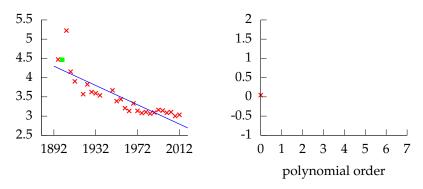




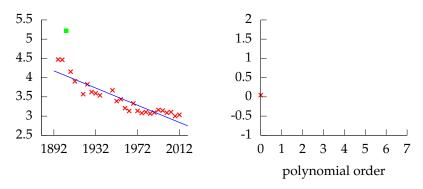




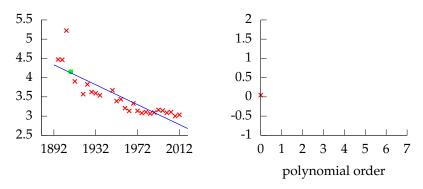
Polynomial order 1, training error -21.183, leave one out error -0.15413.



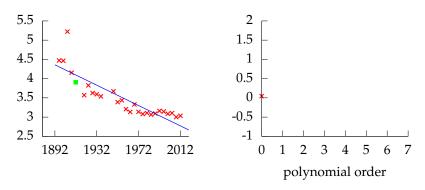
Polynomial order 1, training error -21.183, leave one out error -0.15413.



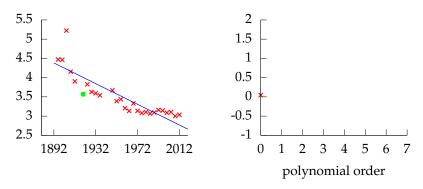
Polynomial order 1, training error -21.183, leave one out error -0.15413.



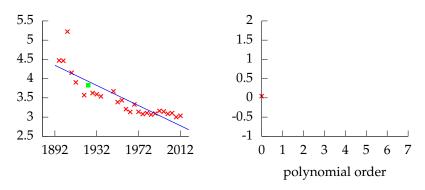
Polynomial order 1, training error -21.183, leave one out error -0.15413.



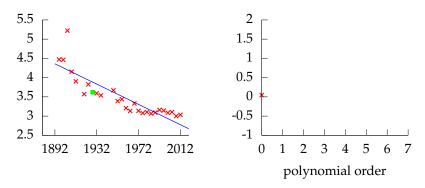
Polynomial order 1, training error -21.183, leave one out error -0.15413.



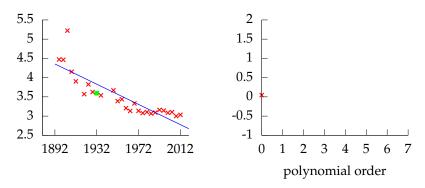
Polynomial order 1, training error -21.183, leave one out error -0.15413.



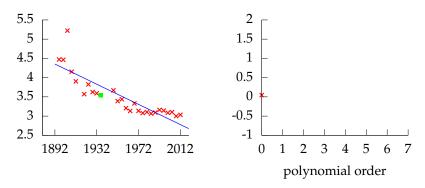
Polynomial order 1, training error -21.183, leave one out error -0.15413.



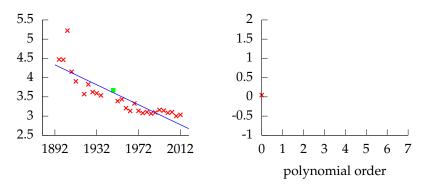
Polynomial order 1, training error -21.183, leave one out error -0.15413.



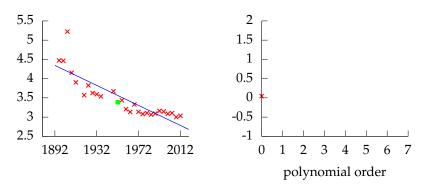
Polynomial order 1, training error -21.183, leave one out error -0.15413.



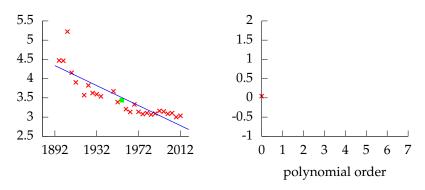
Polynomial order 1, training error -21.183, leave one out error -0.15413.



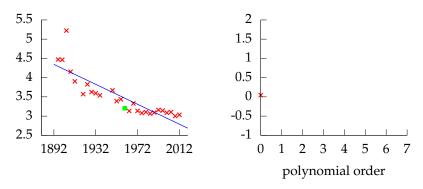
Polynomial order 1, training error -21.183, leave one out error -0.15413.



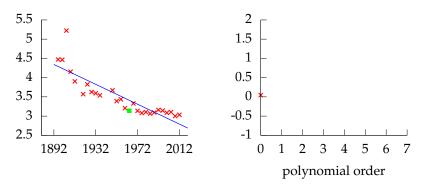
Polynomial order 1, training error -21.183, leave one out error -0.15413.



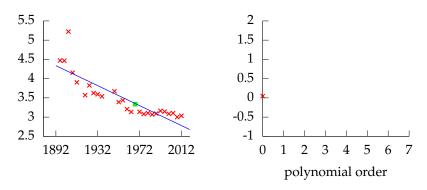
Polynomial order 1, training error -21.183, leave one out error -0.15413.



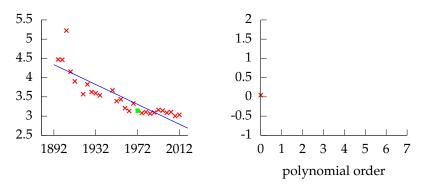
Polynomial order 1, training error -21.183, leave one out error -0.15413.



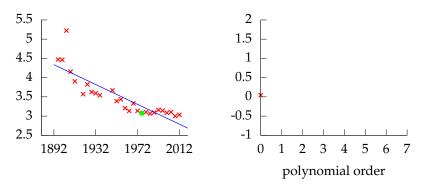
Polynomial order 1, training error -21.183, leave one out error -0.15413.



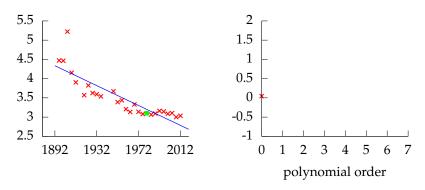
Polynomial order 1, training error -21.183, leave one out error -0.15413.



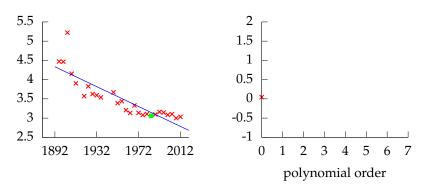
Polynomial order 1, training error -21.183, leave one out error -0.15413.



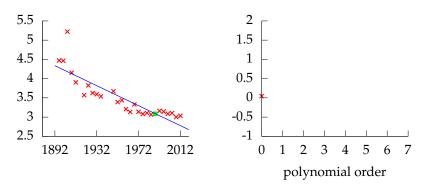
Polynomial order 1, training error -21.183, leave one out error -0.15413.



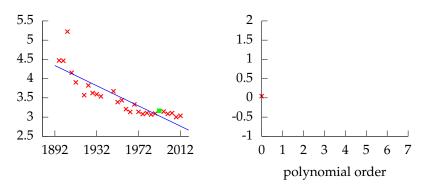
Polynomial order 1, training error -21.183, leave one out error -0.15413.



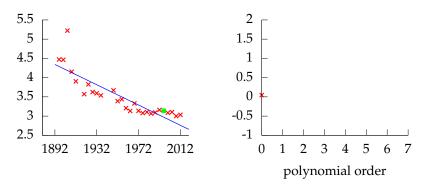
Polynomial order 1, training error -21.183, leave one out error -0.15413.



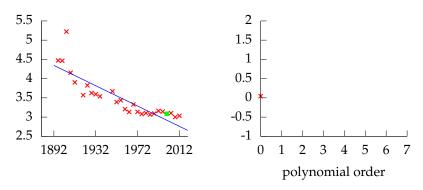
Polynomial order 1, training error -21.183, leave one out error -0.15413.



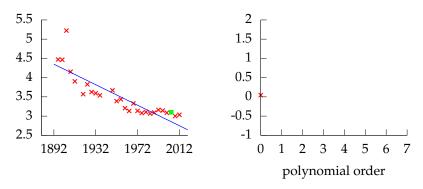
Polynomial order 1, training error -21.183, leave one out error -0.15413.



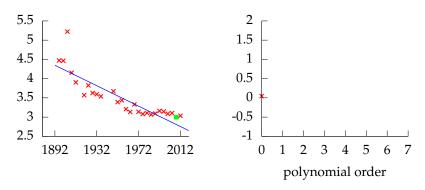
Polynomial order 1, training error -21.183, leave one out error -0.15413.



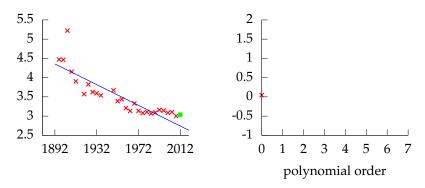
Polynomial order 1, training error -21.183, leave one out error -0.15413.



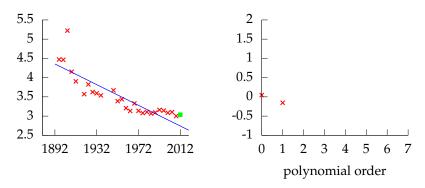
Polynomial order 1, training error -21.183, leave one out error -0.15413.



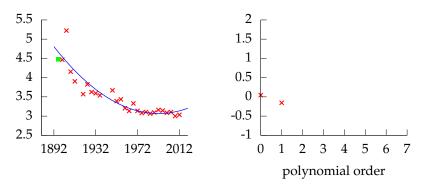
Polynomial order 1, training error -21.183, leave one out error -0.15413.

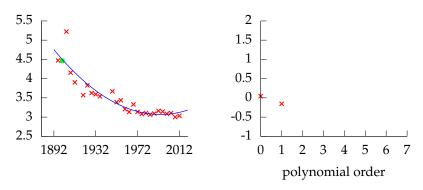


Polynomial order 1, training error -21.183, leave one out error -0.15413.

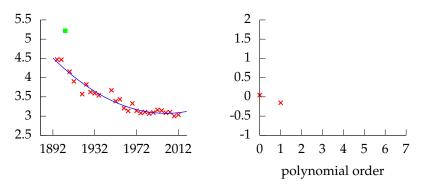


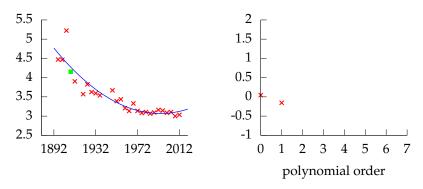
Polynomial order 1, training error -21.183, leave one out error -0.15413.

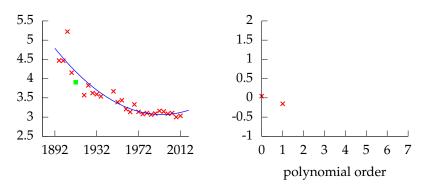




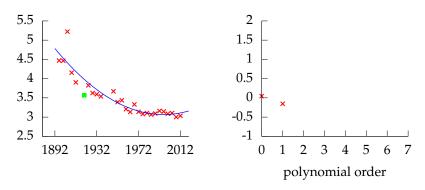
Polynomial order 2, training error -28.403, leave one out error 0.34669.



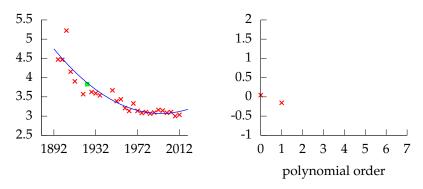


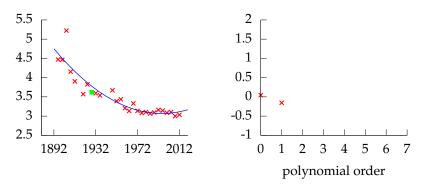


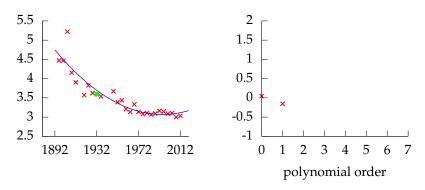
Polynomial order 2, training error -28.403, leave one out error 0.34669.



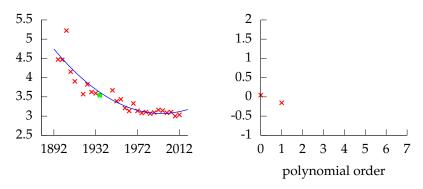
Polynomial order 2, training error -28.403, leave one out error 0.34669.

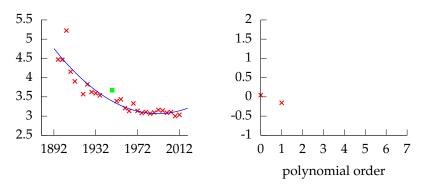


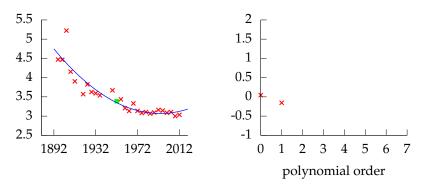


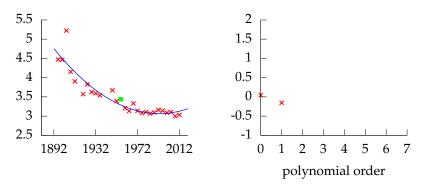


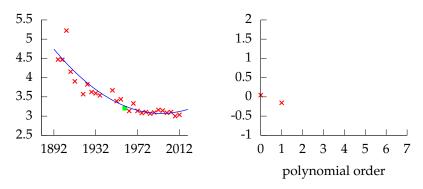
Polynomial order 2, training error -28.403, leave one out error 0.34669.

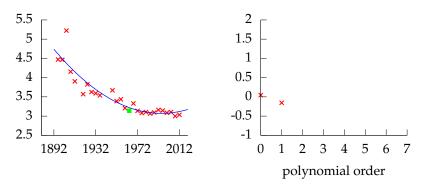


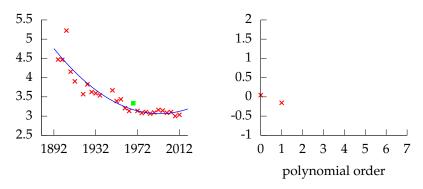


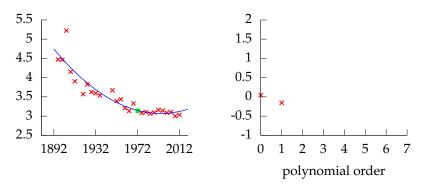


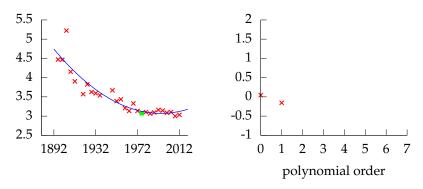


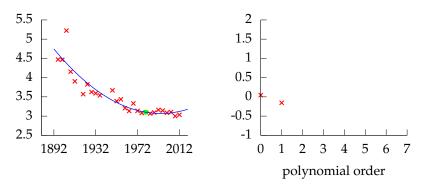


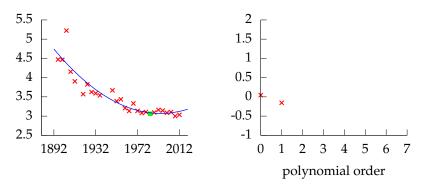


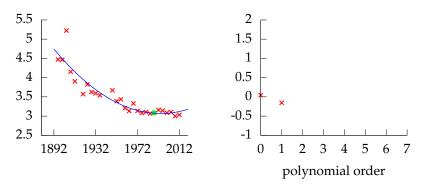


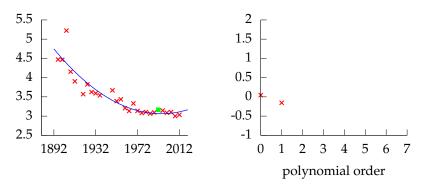


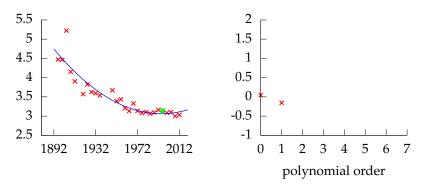


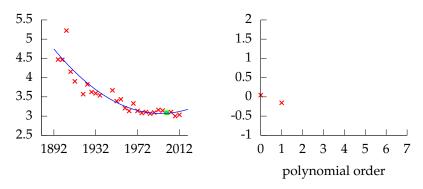


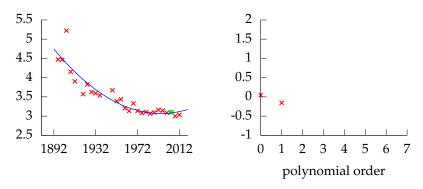


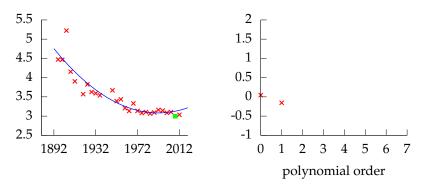


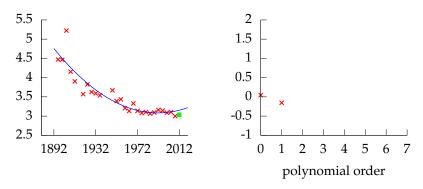


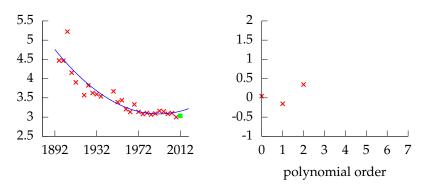




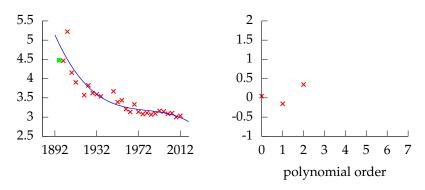




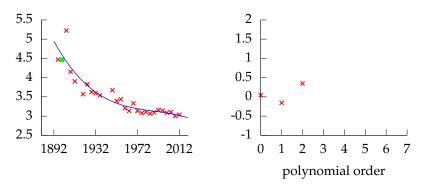




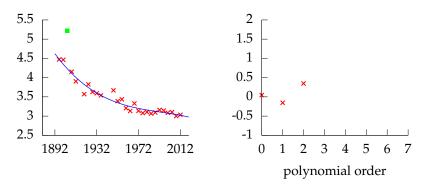
Polynomial order 2, training error -28.403, leave one out error 0.34669.



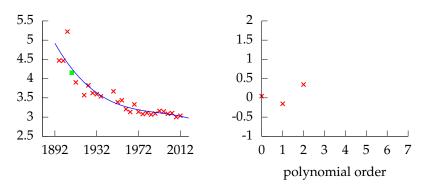
Polynomial order 3, training error -29.223, leave one out error 0.51621.



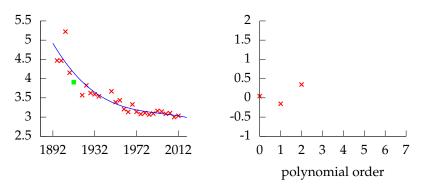
Polynomial order 3, training error -29.223, leave one out error 0.51621.



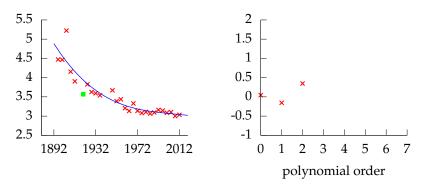
Polynomial order 3, training error -29.223, leave one out error 0.51621.



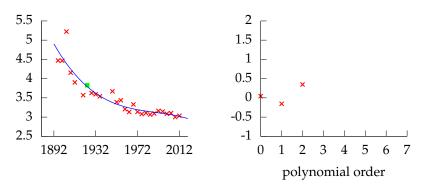
Polynomial order 3, training error -29.223, leave one out error 0.51621.



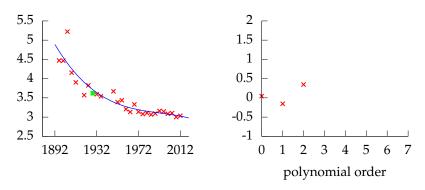
Polynomial order 3, training error -29.223, leave one out error 0.51621.



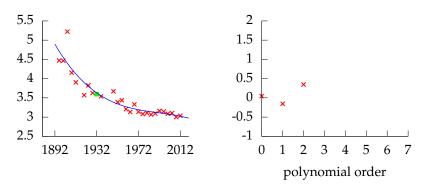
Polynomial order 3, training error -29.223, leave one out error 0.51621.



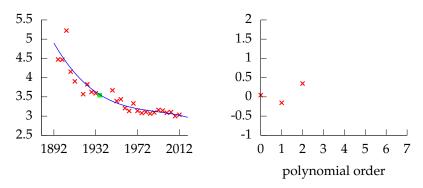
Polynomial order 3, training error -29.223, leave one out error 0.51621.



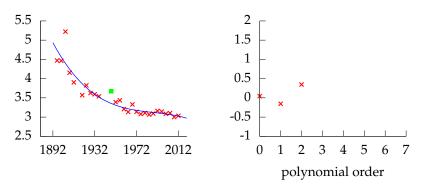
Polynomial order 3, training error -29.223, leave one out error 0.51621.



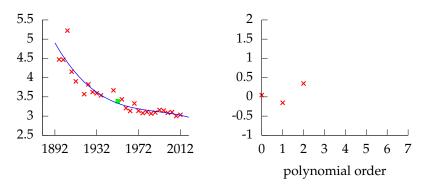
Polynomial order 3, training error -29.223, leave one out error 0.51621.



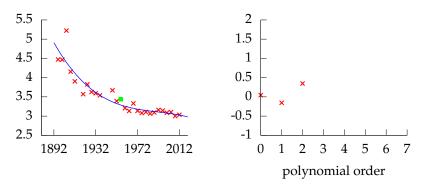
Polynomial order 3, training error -29.223, leave one out error 0.51621.



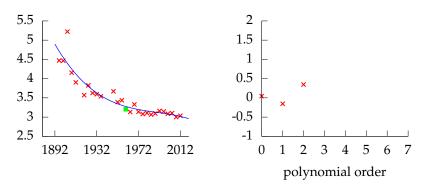
Polynomial order 3, training error -29.223, leave one out error 0.51621.



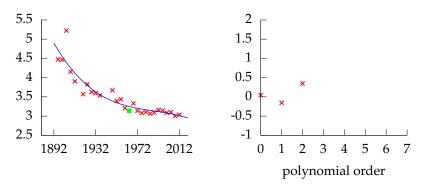
Polynomial order 3, training error -29.223, leave one out error 0.51621.



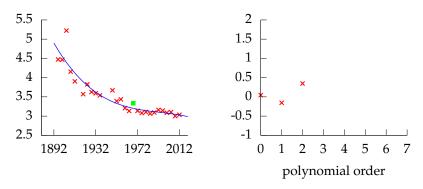
Polynomial order 3, training error -29.223, leave one out error 0.51621.



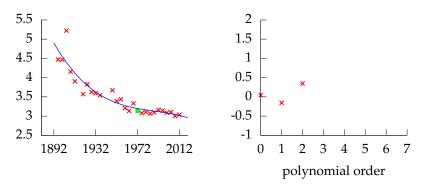
Polynomial order 3, training error -29.223, leave one out error 0.51621.



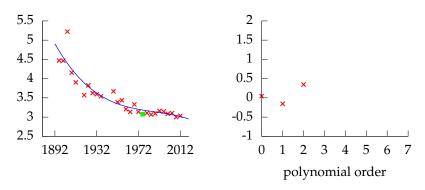
Polynomial order 3, training error -29.223, leave one out error 0.51621.



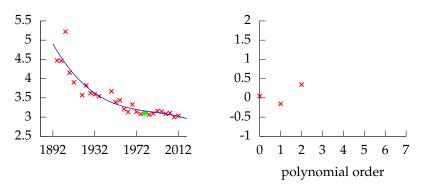
Polynomial order 3, training error -29.223, leave one out error 0.51621.



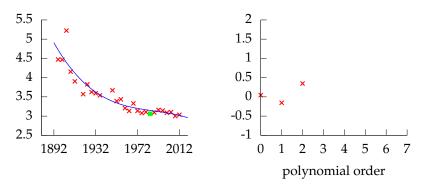
Polynomial order 3, training error -29.223, leave one out error 0.51621.



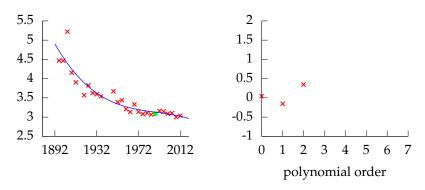
Polynomial order 3, training error -29.223, leave one out error 0.51621.



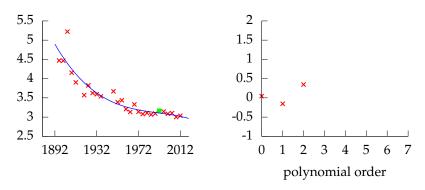
Polynomial order 3, training error -29.223, leave one out error 0.51621.



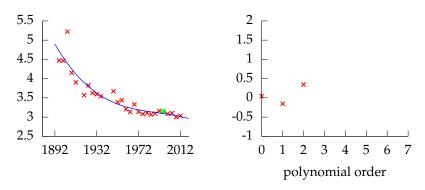
Polynomial order 3, training error -29.223, leave one out error 0.51621.



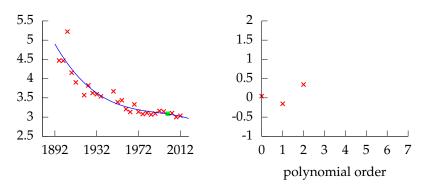
Polynomial order 3, training error -29.223, leave one out error 0.51621.



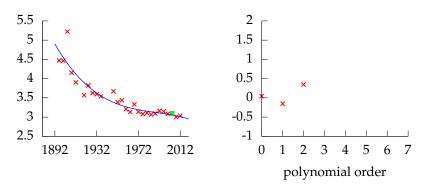
Polynomial order 3, training error -29.223, leave one out error 0.51621.



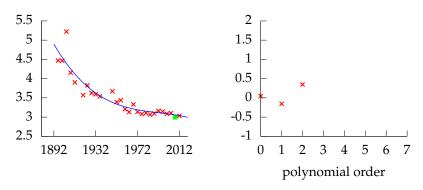
Polynomial order 3, training error -29.223, leave one out error 0.51621.



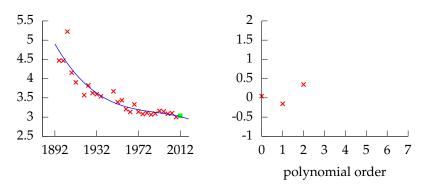
Polynomial order 3, training error -29.223, leave one out error 0.51621.



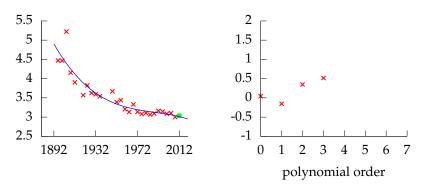
Polynomial order 3, training error -29.223, leave one out error 0.51621.



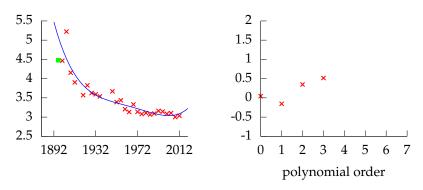
Polynomial order 3, training error -29.223, leave one out error 0.51621.



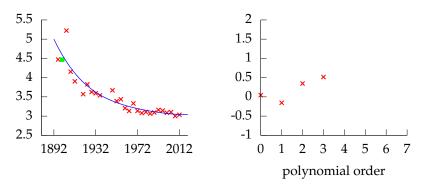
Polynomial order 3, training error -29.223, leave one out error 0.51621.



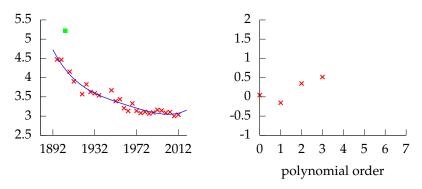
Polynomial order 3, training error -29.223, leave one out error 0.51621.



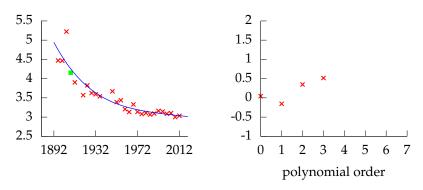
Polynomial order 4, training error -29.324, leave one out error 0.84844.



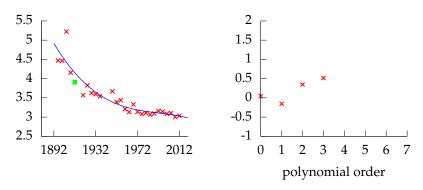
Polynomial order 4, training error -29.324, leave one out error 0.84844.



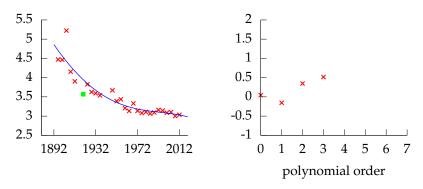
Polynomial order 4, training error -29.324, leave one out error 0.84844.



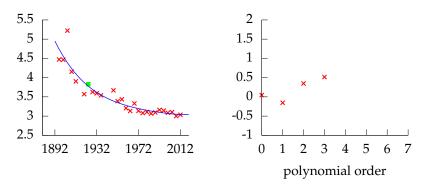
Polynomial order 4, training error -29.324, leave one out error 0.84844.



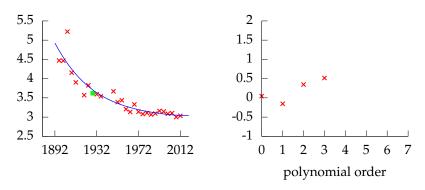
Polynomial order 4, training error -29.324, leave one out error 0.84844.



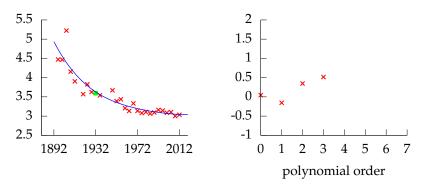
Polynomial order 4, training error -29.324, leave one out error 0.84844.



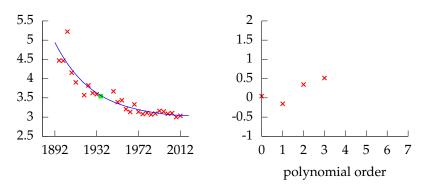
Polynomial order 4, training error -29.324, leave one out error 0.84844.



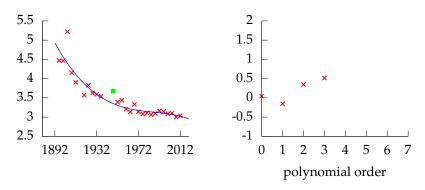
Polynomial order 4, training error -29.324, leave one out error 0.84844.



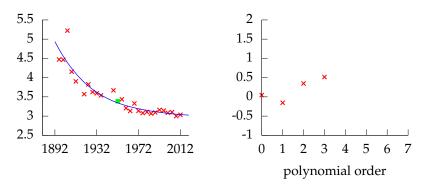
Polynomial order 4, training error -29.324, leave one out error 0.84844.



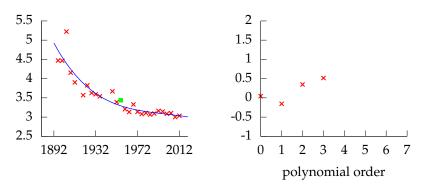
Polynomial order 4, training error -29.324, leave one out error 0.84844.



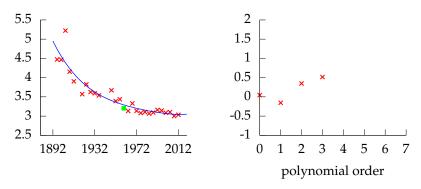
Polynomial order 4, training error -29.324, leave one out error 0.84844.



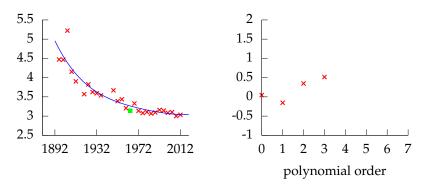
Polynomial order 4, training error -29.324, leave one out error 0.84844.



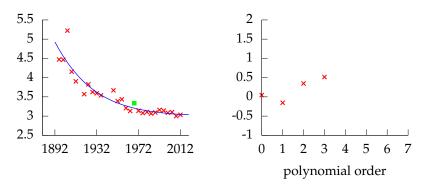
Polynomial order 4, training error -29.324, leave one out error 0.84844.



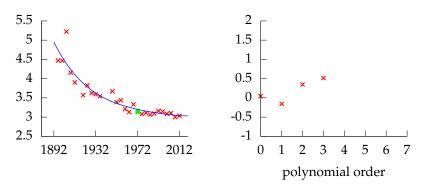
Polynomial order 4, training error -29.324, leave one out error 0.84844.



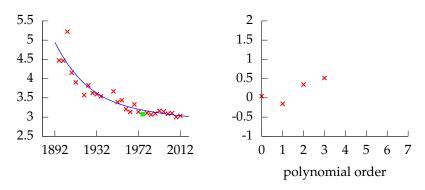
Polynomial order 4, training error -29.324, leave one out error 0.84844.



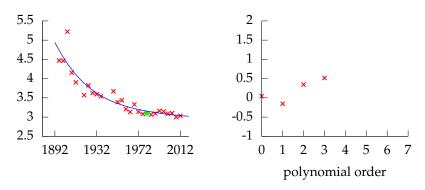
Polynomial order 4, training error -29.324, leave one out error 0.84844.



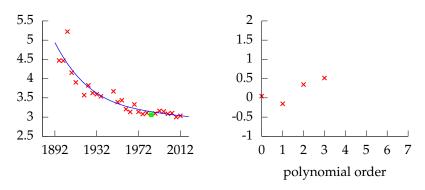
Polynomial order 4, training error -29.324, leave one out error 0.84844.



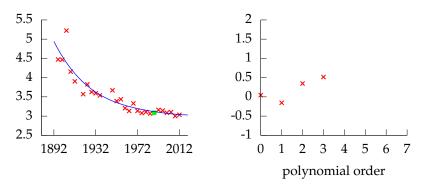
Polynomial order 4, training error -29.324, leave one out error 0.84844.



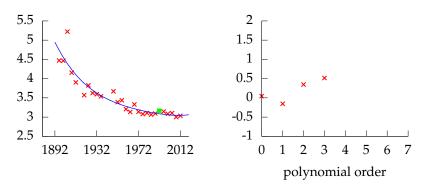
Polynomial order 4, training error -29.324, leave one out error 0.84844.



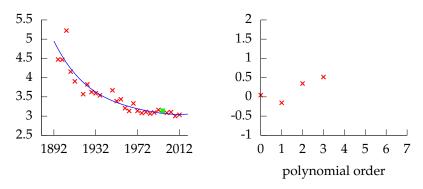
Polynomial order 4, training error -29.324, leave one out error 0.84844.



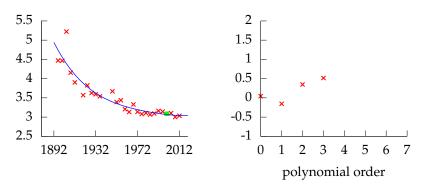
Polynomial order 4, training error -29.324, leave one out error 0.84844.



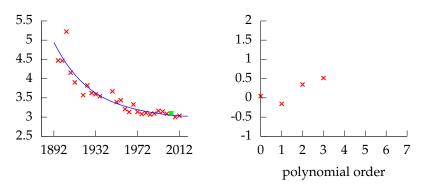
Polynomial order 4, training error -29.324, leave one out error 0.84844.



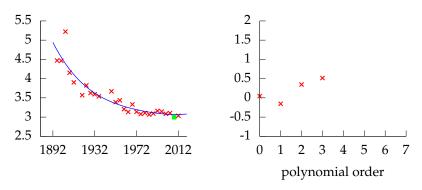
Polynomial order 4, training error -29.324, leave one out error 0.84844.



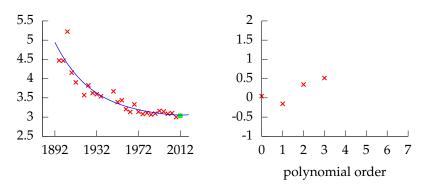
Polynomial order 4, training error -29.324, leave one out error 0.84844.



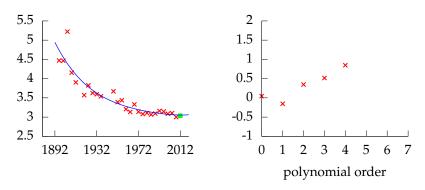
Polynomial order 4, training error -29.324, leave one out error 0.84844.



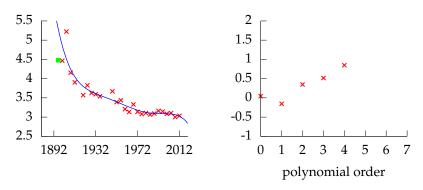
Polynomial order 4, training error -29.324, leave one out error 0.84844.



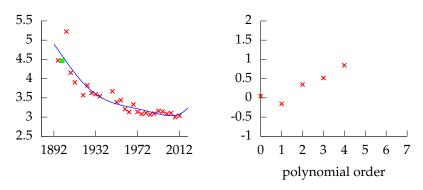
Polynomial order 4, training error -29.324, leave one out error 0.84844.



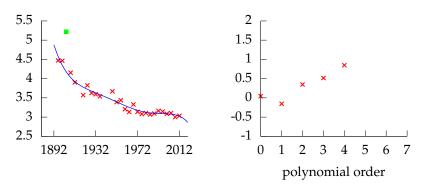
Polynomial order 4, training error -29.324, leave one out error 0.84844.



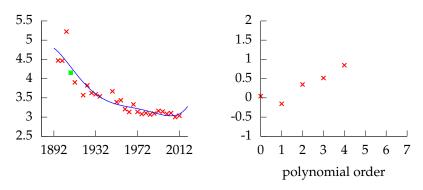
Polynomial order 5, training error -29.524, leave one out error 1.48.



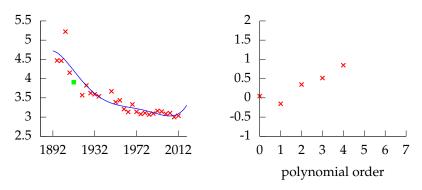
Polynomial order 5, training error -29.524, leave one out error 1.48.



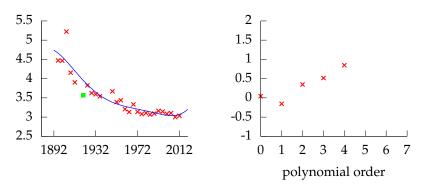
Polynomial order 5, training error -29.524, leave one out error 1.48.



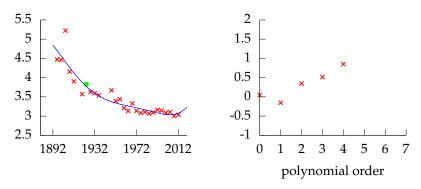
Polynomial order 5, training error -29.524, leave one out error 1.48.



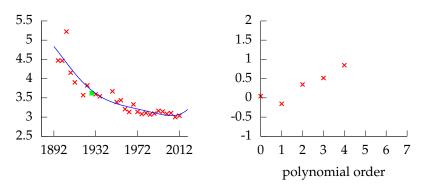
Polynomial order 5, training error -29.524, leave one out error 1.48.



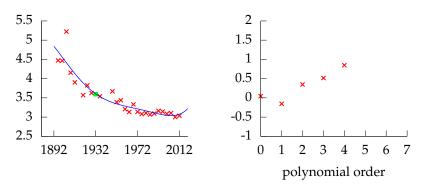
Polynomial order 5, training error -29.524, leave one out error 1.48.



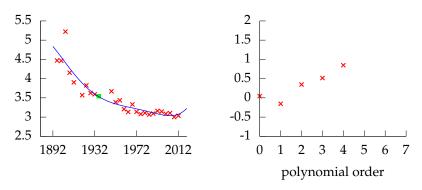
Polynomial order 5, training error -29.524, leave one out error 1.48.



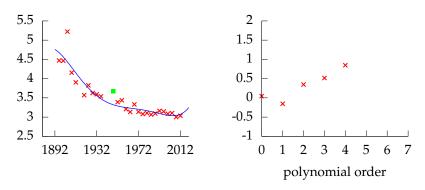
Polynomial order 5, training error -29.524, leave one out error 1.48.



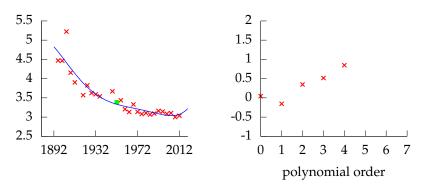
Polynomial order 5, training error -29.524, leave one out error 1.48.



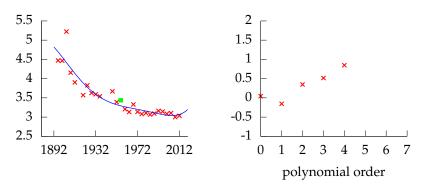
Polynomial order 5, training error -29.524, leave one out error 1.48.



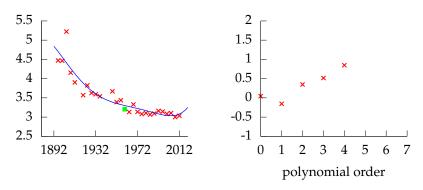
Polynomial order 5, training error -29.524, leave one out error 1.48.



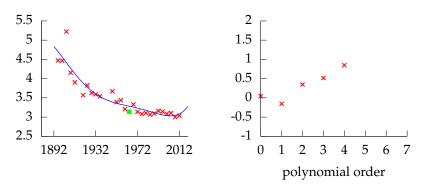
Polynomial order 5, training error -29.524, leave one out error 1.48.



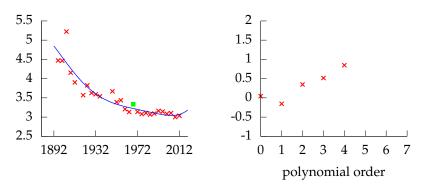
Polynomial order 5, training error -29.524, leave one out error 1.48.



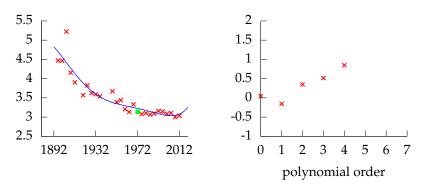
Polynomial order 5, training error -29.524, leave one out error 1.48.



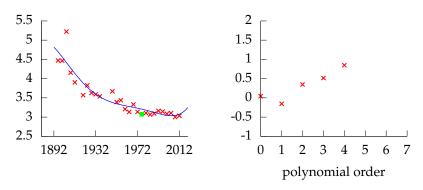
Polynomial order 5, training error -29.524, leave one out error 1.48.



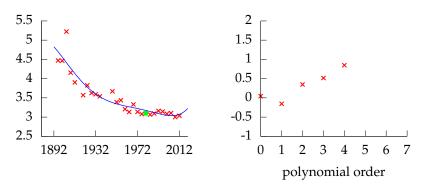
Polynomial order 5, training error -29.524, leave one out error 1.48.



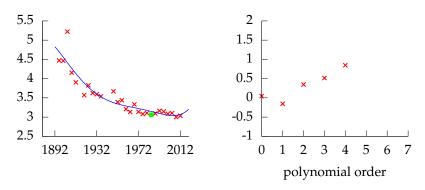
Polynomial order 5, training error -29.524, leave one out error 1.48.



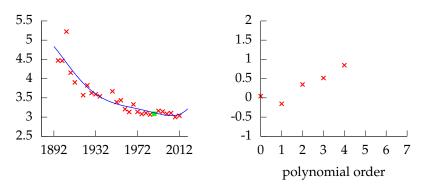
Polynomial order 5, training error -29.524, leave one out error 1.48.



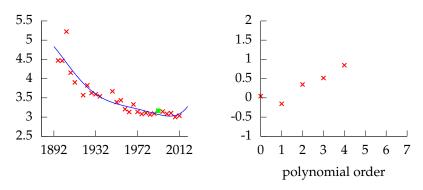
Polynomial order 5, training error -29.524, leave one out error 1.48.



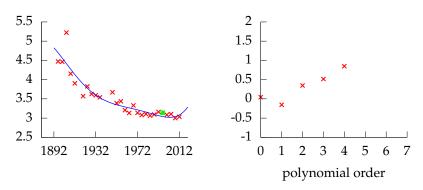
Polynomial order 5, training error -29.524, leave one out error 1.48.



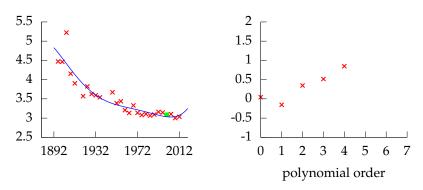
Polynomial order 5, training error -29.524, leave one out error 1.48.



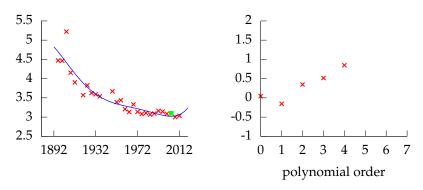
Polynomial order 5, training error -29.524, leave one out error 1.48.



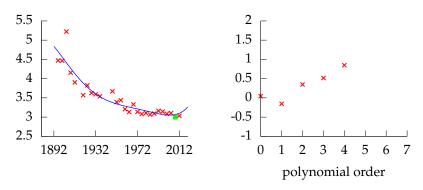
Polynomial order 5, training error -29.524, leave one out error 1.48.



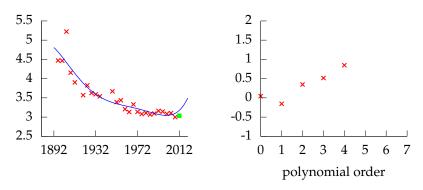
Polynomial order 5, training error -29.524, leave one out error 1.48.



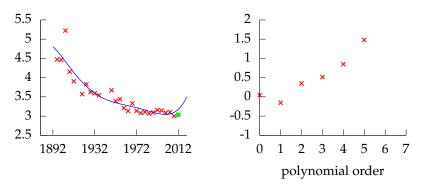
Polynomial order 5, training error -29.524, leave one out error 1.48.



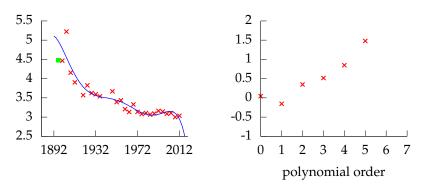
Polynomial order 5, training error -29.524, leave one out error 1.48.



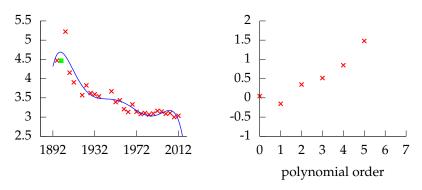
Polynomial order 5, training error -29.524, leave one out error 1.48.



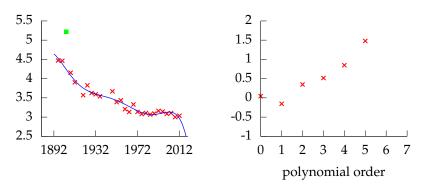
Polynomial order 5, training error -29.524, leave one out error 1.48.



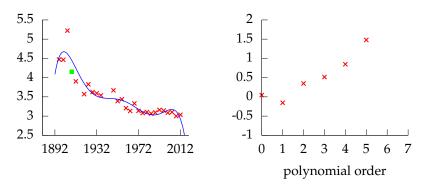
Polynomial order 6, training error -32.237, leave one out error 1.5047.



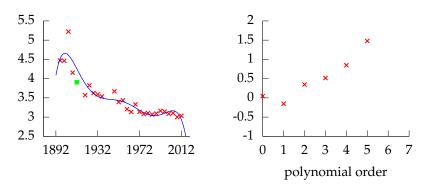
Polynomial order 6, training error -32.237, leave one out error 1.5047.



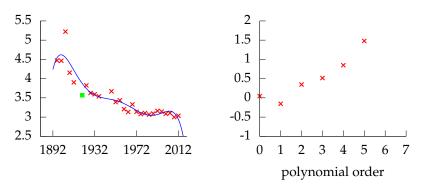
Polynomial order 6, training error -32.237, leave one out error 1.5047.



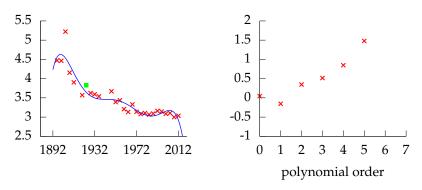
Polynomial order 6, training error -32.237, leave one out error 1.5047.



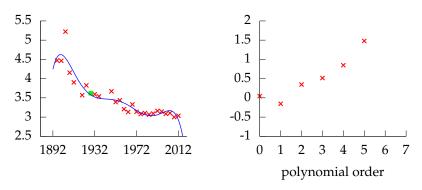
Polynomial order 6, training error -32.237, leave one out error 1.5047.



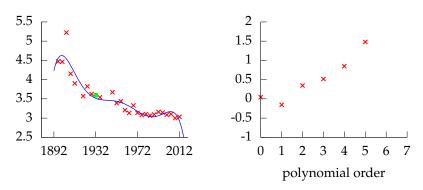
Polynomial order 6, training error -32.237, leave one out error 1.5047.



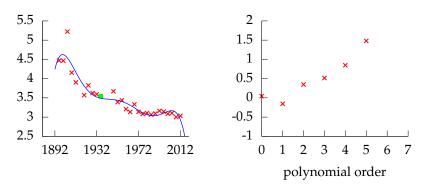
Polynomial order 6, training error -32.237, leave one out error 1.5047.



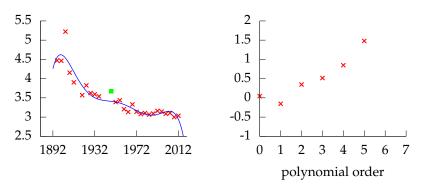
Polynomial order 6, training error -32.237, leave one out error 1.5047.



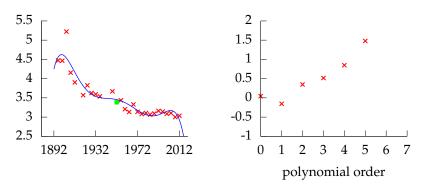
Polynomial order 6, training error -32.237, leave one out error 1.5047.



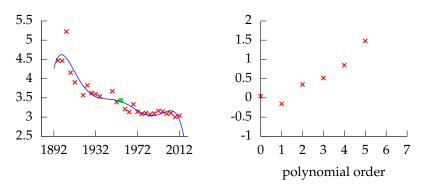
Polynomial order 6, training error -32.237, leave one out error 1.5047.



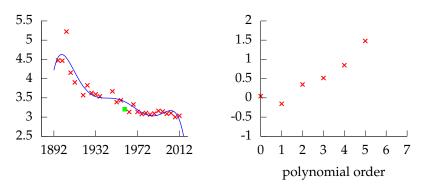
Polynomial order 6, training error -32.237, leave one out error 1.5047.



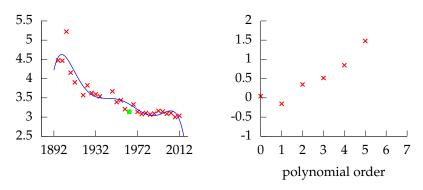
Polynomial order 6, training error -32.237, leave one out error 1.5047.



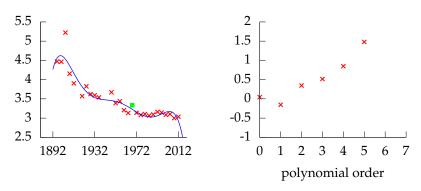
Polynomial order 6, training error -32.237, leave one out error 1.5047.



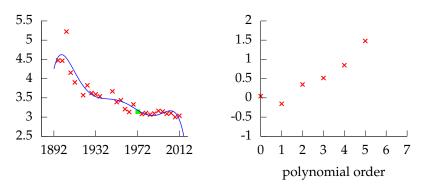
Polynomial order 6, training error -32.237, leave one out error 1.5047.



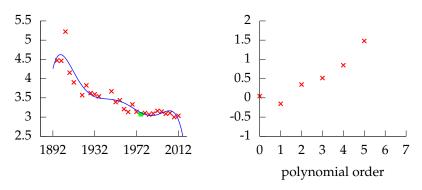
Polynomial order 6, training error -32.237, leave one out error 1.5047.



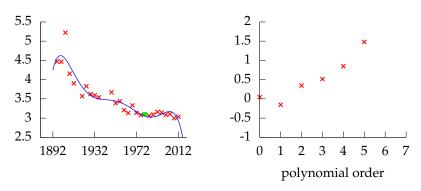
Polynomial order 6, training error -32.237, leave one out error 1.5047.



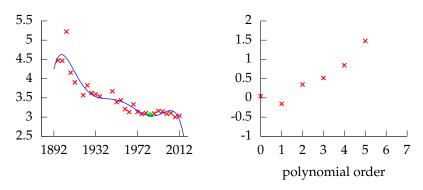
Polynomial order 6, training error -32.237, leave one out error 1.5047.



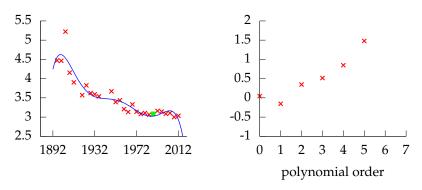
Polynomial order 6, training error -32.237, leave one out error 1.5047.



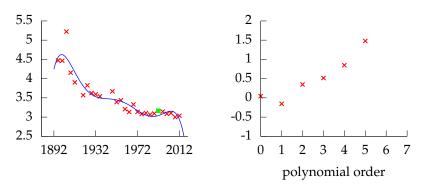
Polynomial order 6, training error -32.237, leave one out error 1.5047.



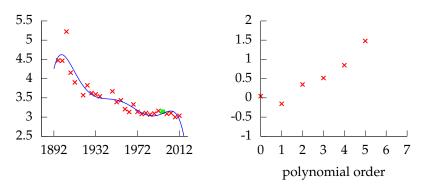
Polynomial order 6, training error -32.237, leave one out error 1.5047.



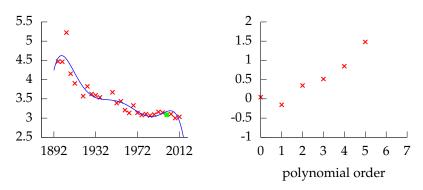
Polynomial order 6, training error -32.237, leave one out error 1.5047.



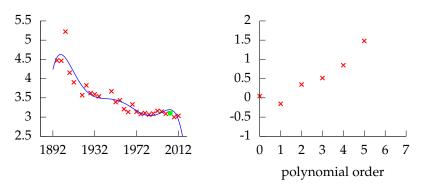
Polynomial order 6, training error -32.237, leave one out error 1.5047.



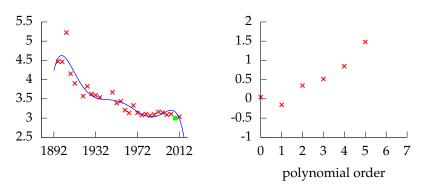
Polynomial order 6, training error -32.237, leave one out error 1.5047.



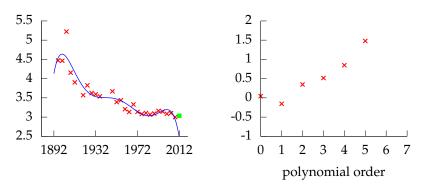
Polynomial order 6, training error -32.237, leave one out error 1.5047.



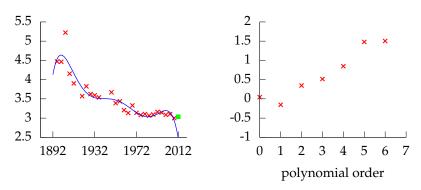
Polynomial order 6, training error -32.237, leave one out error 1.5047.



Polynomial order 6, training error -32.237, leave one out error 1.5047.



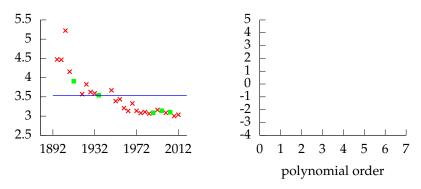
Polynomial order 6, training error -32.237, leave one out error 1.5047.

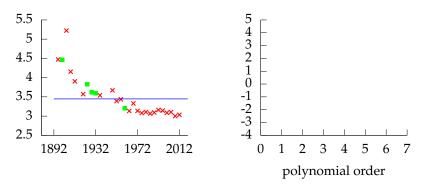


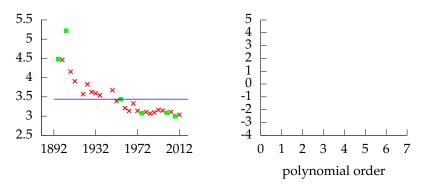
Polynomial order 6, training error -32.237, leave one out error 1.5047.

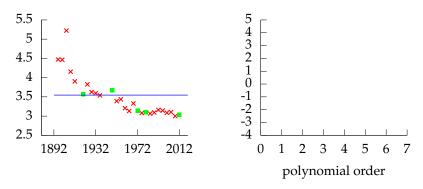
k Fold Cross Validation

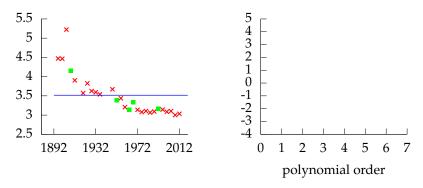
- Leave one out cross validation can be very time consuming!
- ► Need to train your algorithm *n* times.
- ► An alternative: *k* fold cross validation.

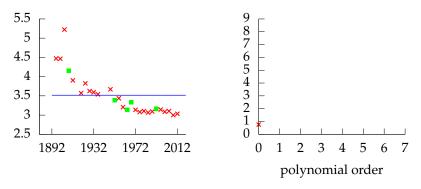


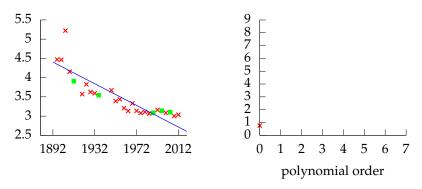




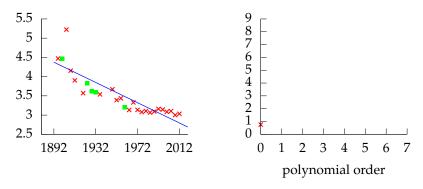




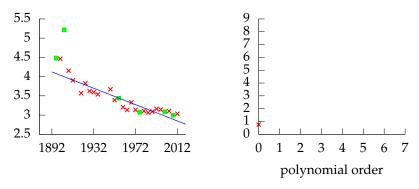




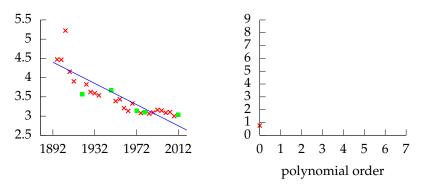
Polynomial order 1, training error -18.873, leave one out error -0.15413.



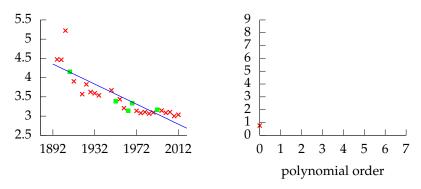
Polynomial order 1, training error -18.873, leave one out error -0.15413.



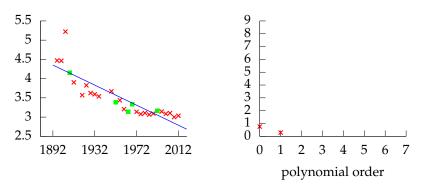
Polynomial order 1, training error -18.873, leave one out error -0.15413.



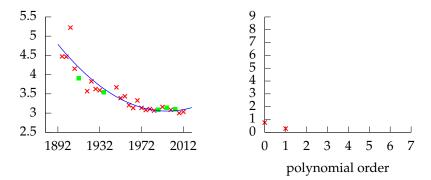
Polynomial order 1, training error -18.873, leave one out error -0.15413.



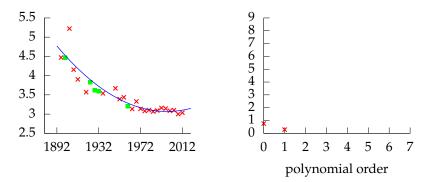
Polynomial order 1, training error -18.873, leave one out error -0.15413.



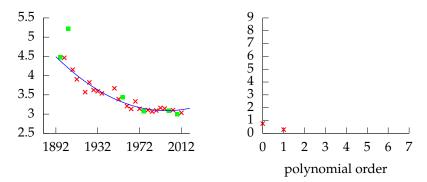
Polynomial order 1, training error -18.873, leave one out error -0.15413.



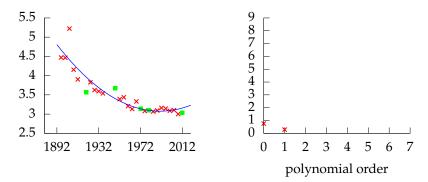
Polynomial order 2, training error -25.177, leave one out error 0.34669.



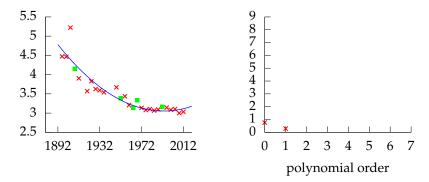
Polynomial order 2, training error -25.177, leave one out error 0.34669.



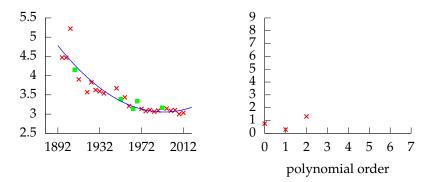
Polynomial order 2, training error -25.177, leave one out error 0.34669.



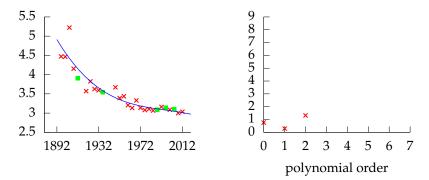
Polynomial order 2, training error -25.177, leave one out error 0.34669.



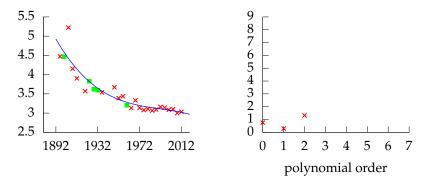
Polynomial order 2, training error -25.177, leave one out error 0.34669.



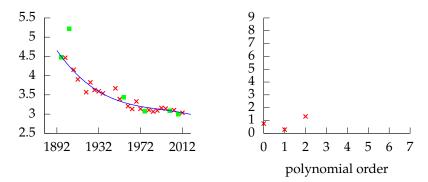
Polynomial order 2, training error -25.177, leave one out error 0.34669.



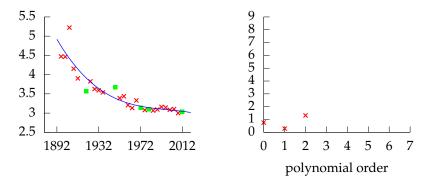
Polynomial order 3, training error -25.777, leave one out error 0.51621.



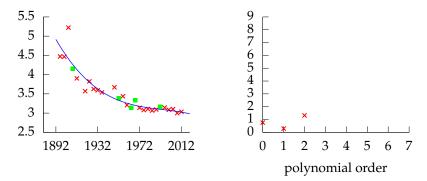
Polynomial order 3, training error -25.777, leave one out error 0.51621.



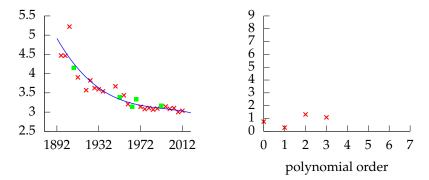
Polynomial order 3, training error -25.777, leave one out error 0.51621.



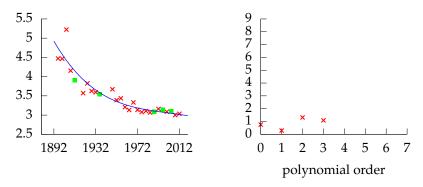
Polynomial order 3, training error -25.777, leave one out error 0.51621.



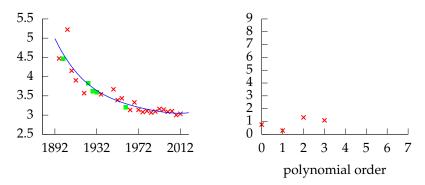
Polynomial order 3, training error -25.777, leave one out error 0.51621.



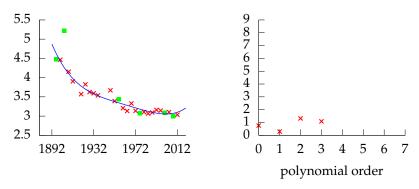
Polynomial order 3, training error -25.777, leave one out error 0.51621.



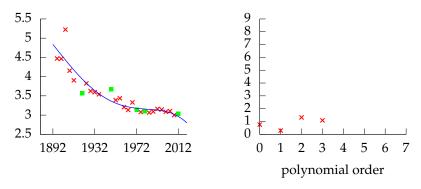
Polynomial order 4, training error -26.048, leave one out error 0.84844.



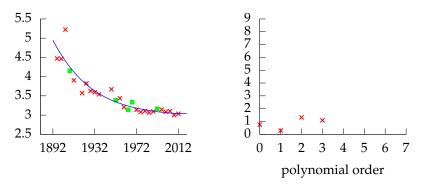
Polynomial order 4, training error -26.048, leave one out error 0.84844.



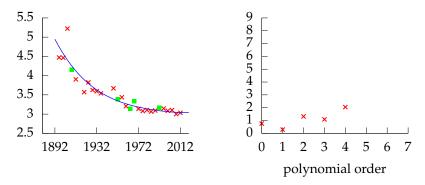
Polynomial order 4, training error -26.048, leave one out error 0.84844.



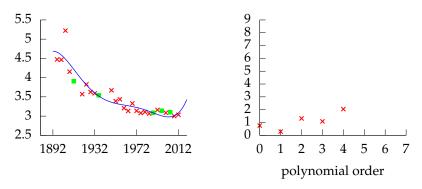
Polynomial order 4, training error -26.048, leave one out error 0.84844.



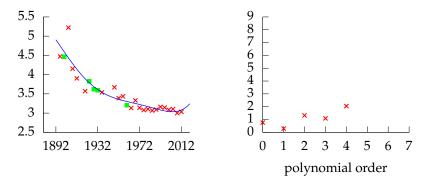
Polynomial order 4, training error -26.048, leave one out error 0.84844.



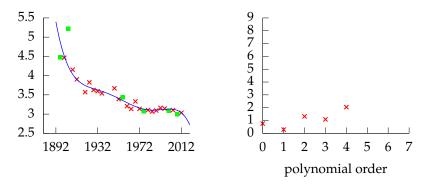
Polynomial order 4, training error -26.048, leave one out error 0.84844.



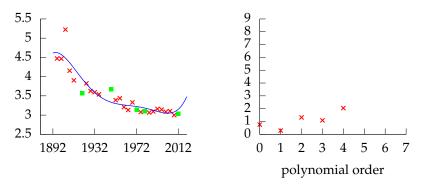
Polynomial order 5, training error -26.892, leave one out error 1.48.



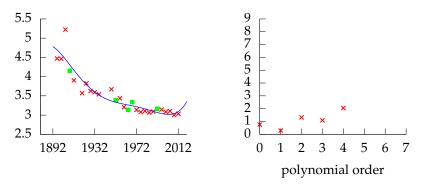
Polynomial order 5, training error -26.892, leave one out error 1.48.



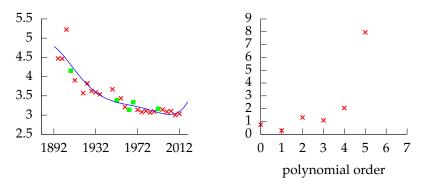
Polynomial order 5, training error -26.892, leave one out error 1.48.



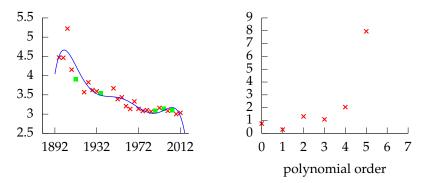
Polynomial order 5, training error -26.892, leave one out error 1.48.



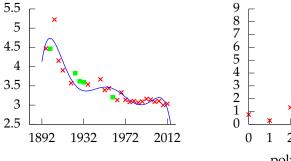
Polynomial order 5, training error -26.892, leave one out error 1.48.

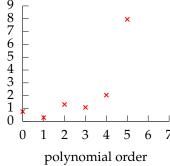


Polynomial order 5, training error -26.892, leave one out error 1.48.

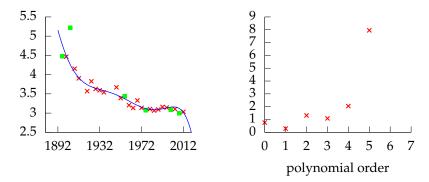


Polynomial order 6, training error -29.395, leave one out error 1.5047.

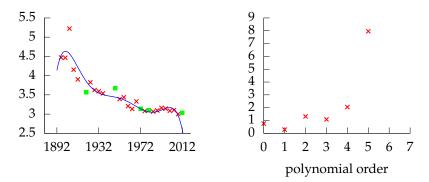




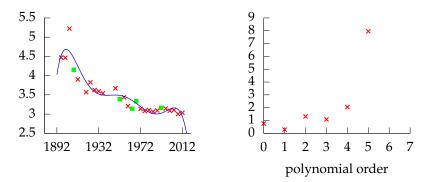
Polynomial order 6, training error -29.395, leave one out error 1.5047.



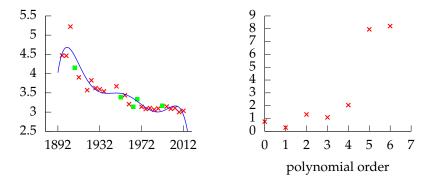
Polynomial order 6, training error -29.395, leave one out error 1.5047.



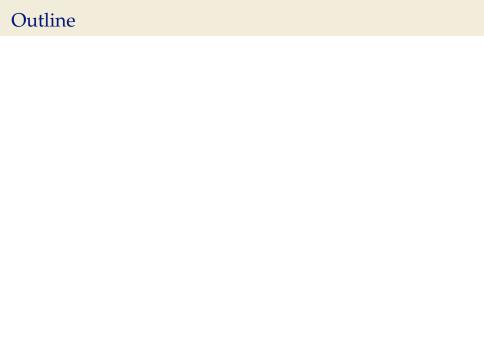
Polynomial order 6, training error -29.395, leave one out error 1.5047.



Polynomial order 6, training error -29.395, leave one out error 1.5047.

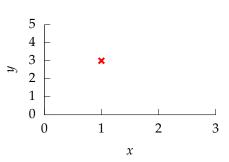


Polynomial order 6, training error -29.395, leave one out error 1.5047.

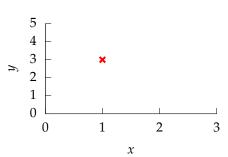


What about two unknowns and *one* observation?

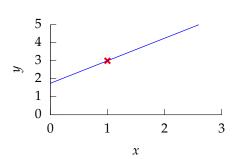
$$y_1 = mx_1 + c$$



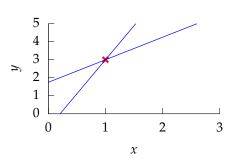
$$m=\frac{y_1-c}{r}$$



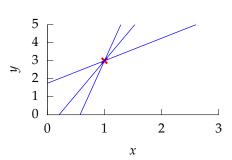
$$c = 1.75 \Longrightarrow m = 1.25$$



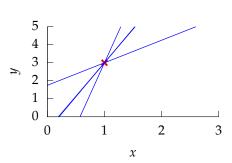
$$c = -0.777 \Longrightarrow m = 3.78$$



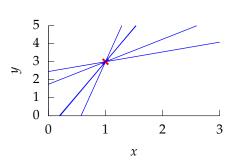
$$c = -4.01 \Longrightarrow m = 7.01$$



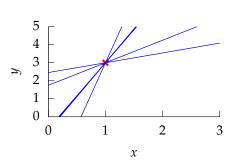
$$c = -0.718 \Longrightarrow m = 3.72$$



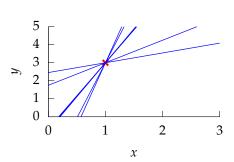
$$c = 2.45 \Longrightarrow m = 0.545$$



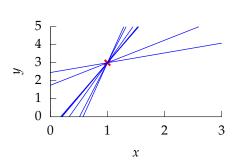
$$c = -0.657 \Longrightarrow m = 3.66$$



$$c = -3.13 \Longrightarrow m = 6.13$$



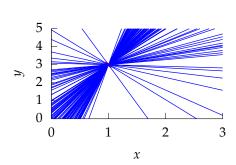
$$c = -1.47 \Longrightarrow m = 4.47$$



Can compute m given c. Assume

$$c \sim \mathcal{N}(0,4)$$
,

we find a distribution of solutions.

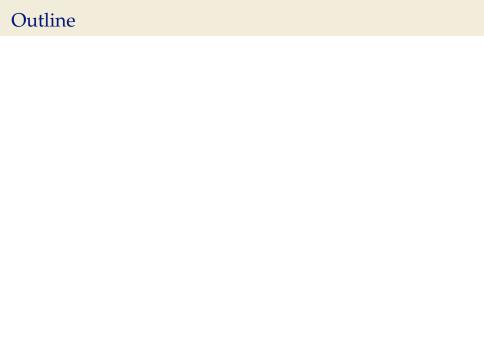


Probability for Under- and Overdetermined

- ▶ To deal with overdetermined introduced probability distribution for 'variable', ϵ_i .
- ► For underdetermined system introduced probability distribution for 'parameter', *c*.
- ► This is known as a Bayesian treatment.

Reading

- ▶ Bishop Section 1.2.3 (pg 21–24).
- ► Bishop Section 1.2.6 (start from just past eq 1.64 pg 30-32).
- ▶ Rogers and Girolami use an example of a coin toss for introducing Bayesian inference Chapter 3, Sections 3.1-3.4 (pg 95-117). Although you also need the beta density which we haven't yet discussed. This is also the example that Laplace used.



Prior Distribution

- ▶ Bayesian inference requires a prior on the parameters.
- ► The prior represents your belief *before* you see the data of the likely value of the parameters.
- ► For linear regression, consider a Gaussian prior on the intercept:

$$c \sim \mathcal{N}(0, \alpha_1)$$

Posterior Distribution

- ► Posterior distribution is found by combining the prior with the likelihood.
- Posterior distribution is your belief after you see the data of the likely value of the parameters.
- ► The posterior is found through Bayes' Rule

$$p(c|y) = \frac{p(y|c)p(c)}{p(y)}$$

Bayes Update

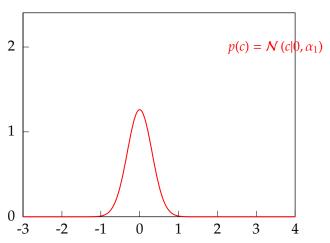


Figure : A Gaussian prior combines with a Gaussian likelihood for a Gaussian posterior.

Bayes Update

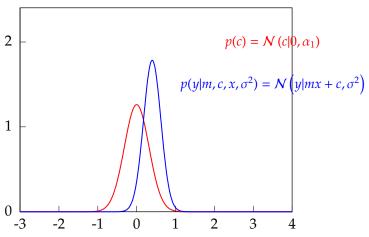


Figure : A Gaussian prior combines with a Gaussian likelihood for a Gaussian posterior.

Bayes Update

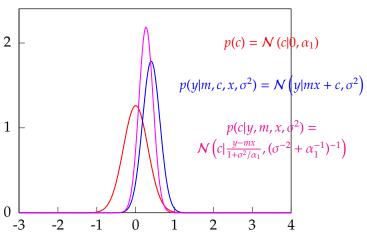


Figure : A Gaussian prior combines with a Gaussian likelihood for a Gaussian posterior.

Stages to Derivation of the Posterior

- Multiply likelihood by prior
 - they are "exponentiated quadratics", the answer is always also an exponentiated quadratic because $\exp(a^2) \exp(b^2) = \exp(a^2 + b^2)$.
- Complete the square to get the resulting density in the form of a Gaussian.
- Recognise the mean and (co)variance of the Gaussian. This is the estimate of the posterior.

Multivariate Prior Distributions

- ► For general Bayesian inference need multivariate priors.
- E.g. for multivariate linear regression:

$$y_i = \sum_i w_j x_{i,j} + \epsilon_i$$

(where we've dropped *c* for convenience), we need a prior over **w**.

- ► This motivates a *multivariate* Gaussian density.
- ▶ We will use the multivariate Gaussian to put a prior *directly* on the function (a Gaussian process).

Multivariate Prior Distributions

- ► For general Bayesian inference need multivariate priors.
- ► E.g. for multivariate linear regression:

$$y_i = \mathbf{w}^{\mathsf{T}} \mathbf{x}_{i,:} + \epsilon_i$$

(where we've dropped *c* for convenience), we need a prior over **w**.

- ► This motivates a *multivariate* Gaussian density.
- ▶ We will use the multivariate Gaussian to put a prior *directly* on the function (a Gaussian process).

Two Dimensional Gaussian

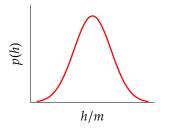
- ▶ Consider height, h/m and weight, w/kg.
- ► Could sample height from a distribution:

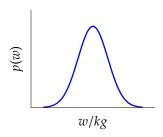
$$p(h) \sim \mathcal{N}(1.7, 0.0225)$$

► And similarly weight:

$$p(w) \sim \mathcal{N}(75, 36)$$

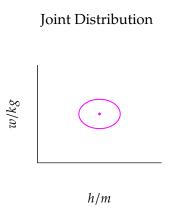
Height and Weight Models

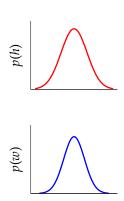




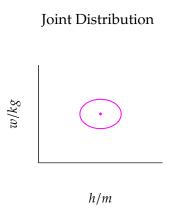
Gaussian distributions for height and weight.

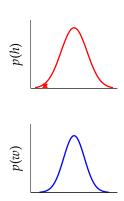
Marginal Distributions



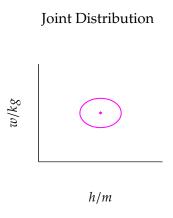


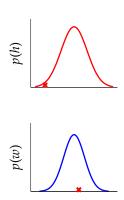
Marginal Distributions



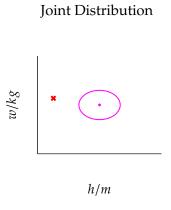


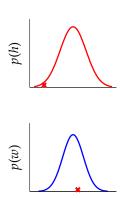
Marginal Distributions



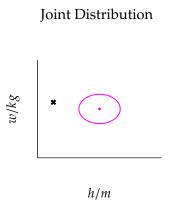


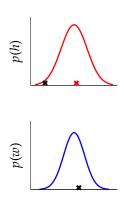
Marginal Distributions



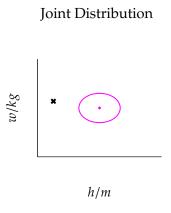


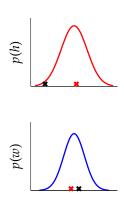
Marginal Distributions



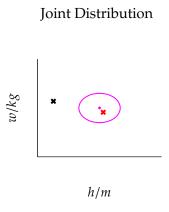


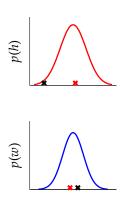
Marginal Distributions



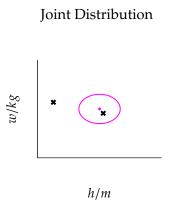


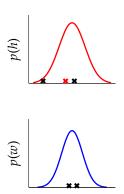
Marginal Distributions



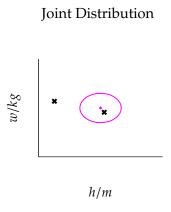


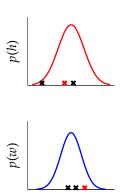
Marginal Distributions



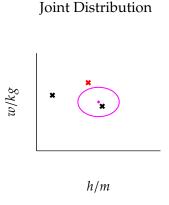


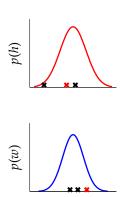
Marginal Distributions



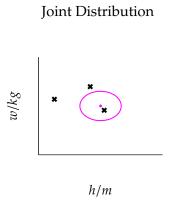


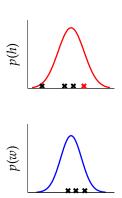
Marginal Distributions



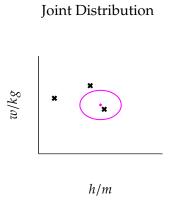


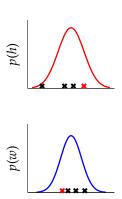
Marginal Distributions



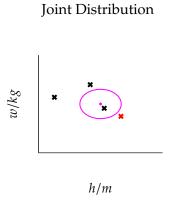


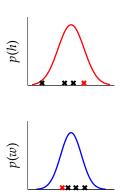
Marginal Distributions



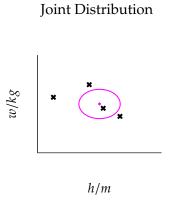


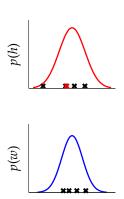
Marginal Distributions



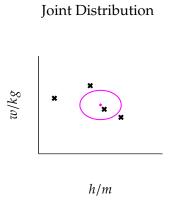


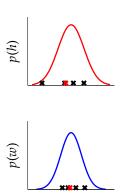
Marginal Distributions



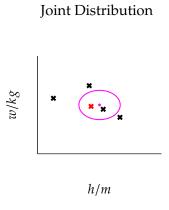


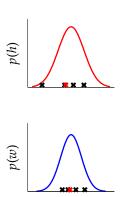
Marginal Distributions



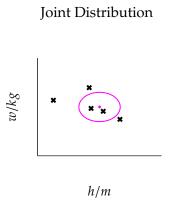


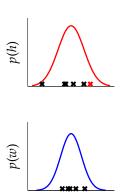
Marginal Distributions



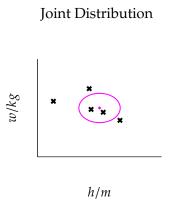


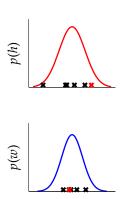
Marginal Distributions



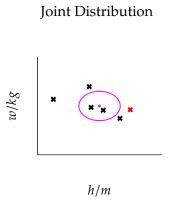


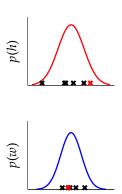
Marginal Distributions



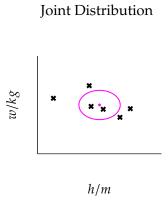


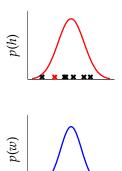
Marginal Distributions



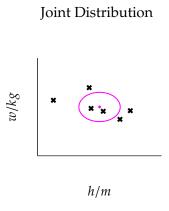


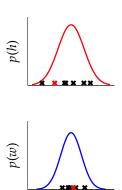
Marginal Distributions



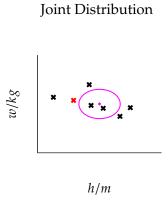


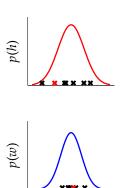
Marginal Distributions



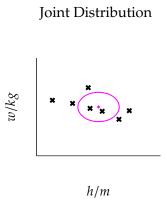


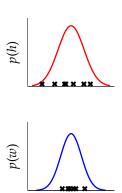
Marginal Distributions





Marginal Distributions



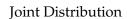


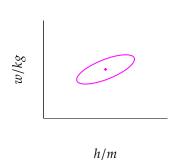
Independence Assumption

► This assumes height and weight are independent.

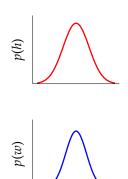
$$p(h, w) = p(h)p(w)$$

► In reality they are dependent (body mass index) = $\frac{w}{h^2}$.

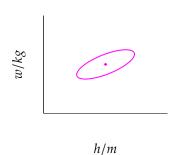




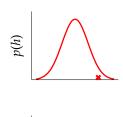
Marginal Distributions

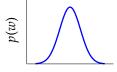


Joint Distribution

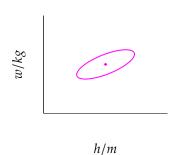


Marginal Distributions

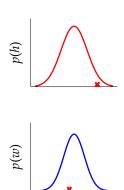




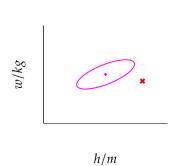


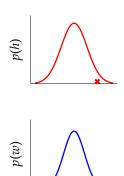


Marginal Distributions

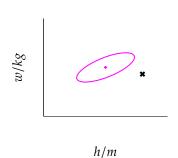


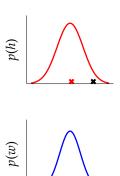
Joint Distribution



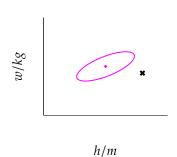


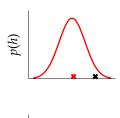
Joint Distribution

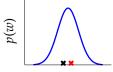




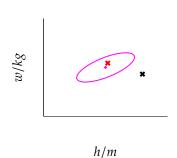
Joint Distribution

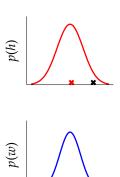




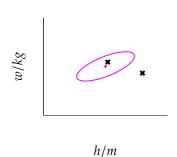


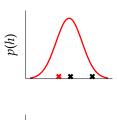
Joint Distribution

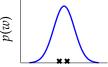




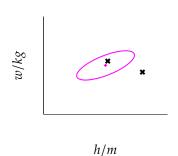
Joint Distribution

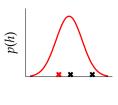


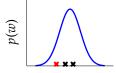




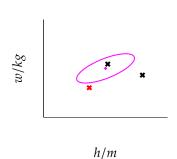
Joint Distribution

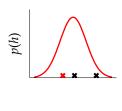


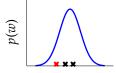




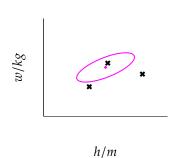
Joint Distribution

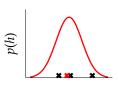


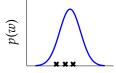




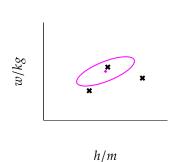
Joint Distribution

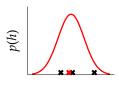


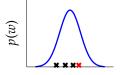




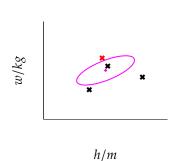
Joint Distribution

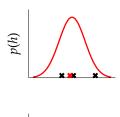


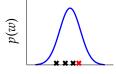




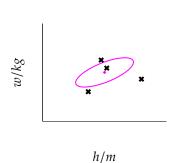
Joint Distribution

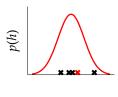


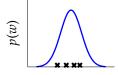




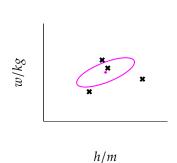
Joint Distribution

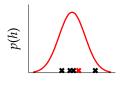


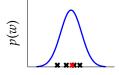




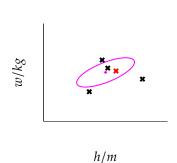
Joint Distribution

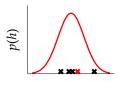


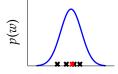




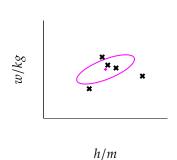
Joint Distribution

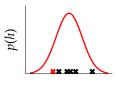


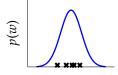




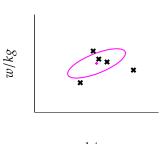
Joint Distribution



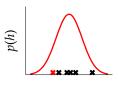


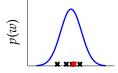


Joint Distribution

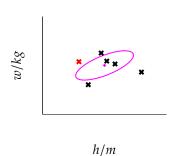


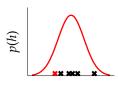
h/m

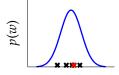




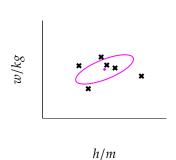
Joint Distribution

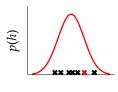


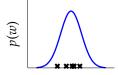




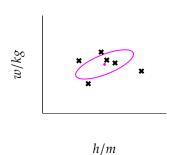
Joint Distribution

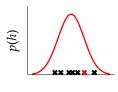


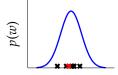




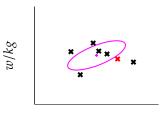
Joint Distribution



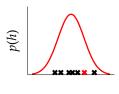


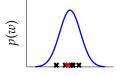


Joint Distribution

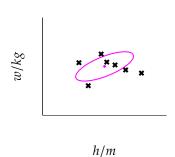


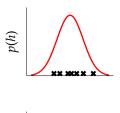
h/m

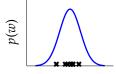




Joint Distribution







$$p(w,h) = p(w)p(h)$$

$$p(w,h) = \frac{1}{\sqrt{2\pi\sigma_1^2}} \left(\frac{1}{\sqrt{2\pi\sigma_2^2}} \exp\left(-\frac{1}{2} \left(\frac{(w-\mu_1)^2}{\sigma_1^2} + \frac{(h-\mu_2)^2}{\sigma_2^2} \right) \right) \right)$$

$$p(w,h) = \frac{1}{2\pi\sqrt{\sigma_1^2\sigma_2^2}} \exp\left(-\frac{1}{2} \begin{pmatrix} \begin{bmatrix} w \\ h \end{bmatrix} - \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix} \end{pmatrix}^{\mathsf{T}} \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix}^{-1} \begin{pmatrix} \begin{bmatrix} w \\ h \end{bmatrix} - \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix} \end{pmatrix}\right)$$

$$p(\mathbf{y}) = \frac{1}{2\pi |\mathbf{D}|} \exp\left(-\frac{1}{2}(\mathbf{y} - \boldsymbol{\mu})^{\top} \mathbf{D}^{-1}(\mathbf{y} - \boldsymbol{\mu})\right)$$

Form correlated from original by rotating the data space using matrix \mathbf{R} .

$$p(\mathbf{y}) = \frac{1}{2\pi |\mathbf{D}|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(\mathbf{y} - \boldsymbol{\mu})^{\mathsf{T}} \mathbf{D}^{-1}(\mathbf{y} - \boldsymbol{\mu})\right)$$

Form correlated from original by rotating the data space using matrix **R**.

$$p(\mathbf{y}) = \frac{1}{2\pi |\mathbf{D}|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(\mathbf{R}^{\top}\mathbf{y} - \mathbf{R}^{\top}\boldsymbol{\mu})^{\top}\mathbf{D}^{-1}(\mathbf{R}^{\top}\mathbf{y} - \mathbf{R}^{\top}\boldsymbol{\mu})\right)$$

Form correlated from original by rotating the data space using matrix \mathbf{R} .

$$p(\mathbf{y}) = \frac{1}{2\pi |\mathbf{D}|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(\mathbf{y} - \boldsymbol{\mu})^{\top} \mathbf{R} \mathbf{D}^{-1} \mathbf{R}^{\top} (\mathbf{y} - \boldsymbol{\mu})\right)$$

this gives a covariance matrix:

$$\mathbf{C}^{-1} = \mathbf{R}\mathbf{D}^{-1}\mathbf{R}^{\mathsf{T}}$$

Form correlated from original by rotating the data space using matrix **R**.

$$p(\mathbf{y}) = \frac{1}{2\pi |\mathbf{C}|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(\mathbf{y} - \boldsymbol{\mu})^{\top} \mathbf{C}^{-1}(\mathbf{y} - \boldsymbol{\mu})\right)$$

this gives a covariance matrix:

$$C = RDR^{T}$$

1. Sum of Gaussian variables is also Gaussian.

$$y_i \sim \mathcal{N}\left(\mu_i, \sigma_i^2\right)$$

1. Sum of Gaussian variables is also Gaussian.

$$y_i \sim \mathcal{N}\left(\mu_i, \sigma_i^2\right)$$

$$\sum_{i=1}^{n} y_i \sim \mathcal{N}\left(\sum_{i=1}^{n} \mu_i, \sum_{i=1}^{n} \sigma_i^2\right)$$

1. Sum of Gaussian variables is also Gaussian.

$$y_i \sim \mathcal{N}\left(\mu_i, \sigma_i^2\right)$$

$$\sum_{i=1}^{n} y_i \sim \mathcal{N}\left(\sum_{i=1}^{n} \mu_i, \sum_{i=1}^{n} \sigma_i^2\right)$$

2. Scaling a Gaussian leads to a Gaussian.

1. Sum of Gaussian variables is also Gaussian.

$$y_i \sim \mathcal{N}\left(\mu_i, \sigma_i^2\right)$$

$$\sum_{i=1}^{n} y_i \sim \mathcal{N}\left(\sum_{i=1}^{n} \mu_i, \sum_{i=1}^{n} \sigma_i^2\right)$$

2. Scaling a Gaussian leads to a Gaussian.

$$y \sim \mathcal{N}\left(\mu, \sigma^2\right)$$

1. Sum of Gaussian variables is also Gaussian.

$$y_i \sim \mathcal{N}\left(\mu_i, \sigma_i^2\right)$$

$$\sum_{i=1}^{n} y_i \sim \mathcal{N}\left(\sum_{i=1}^{n} \mu_i, \sum_{i=1}^{n} \sigma_i^2\right)$$

2. Scaling a Gaussian leads to a Gaussian.

$$y \sim \mathcal{N}\left(\mu, \sigma^2\right)$$

$$wy \sim \mathcal{N}\left(w\mu, w^2\sigma^2\right)$$

Multivariate Consequence

$$\mathbf{x} \sim \mathcal{N}\left(\mu, \mathbf{\Sigma}\right)$$

Multivariate Consequence

$$\mathbf{x} \sim \mathcal{N}\left(\mu, \Sigma\right)$$

► And

$$y = Wx$$

Multivariate Consequence

$$\mathbf{x} \sim \mathcal{N}\left(\mu, \Sigma\right)$$

► And

$$y = Wx$$

► Then

$$\mathbf{y} \sim \mathcal{N}\left(\mathbf{W}\boldsymbol{\mu}, \mathbf{W}\boldsymbol{\Sigma}\mathbf{W}^{\top}\right)$$

References I

- C. M. Bishop. *Pattern Recognition and Machine Learning*. Springer-Verlag, 2006. [Google Books].
- P. S. Laplace. Mémoire sur la probabilité des causes par les évènemens. In *Mémoires de mathèmatique et de physique, presentés à lAcadémie Royale des Sciences, par divers savans, & lù dans ses assemblées 6,* pages 621–656, 1774. Translated in Stigler (1986).
- S. Rogers and M. Girolami. *A First Course in Machine Learning*. CRC Press, 2011. [Google Books].
- S. M. Stigler. Laplace's 1774 memoir on inverse probability. *Statistical Science*, 1:359–378, 1986.