

Bayesian Regression

MLAI: Week 7

Neil D. Lawrence

Department of Computer Science
Sheffield University

11th November 2014

Outline

Quick Review: Overdetermined Systems

Underdetermined Systems

Bayesian Regression

Univariate Bayesian Linear Regression

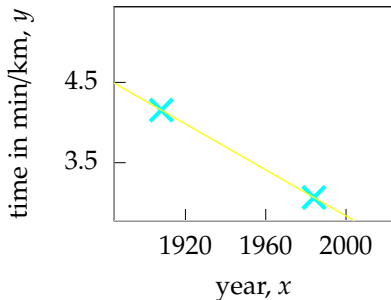
Bayesian Polynomials

Two Simultaneous Equations

A system of two simultaneous equations with two unknowns.

$$y_1 = mx_1 + c$$

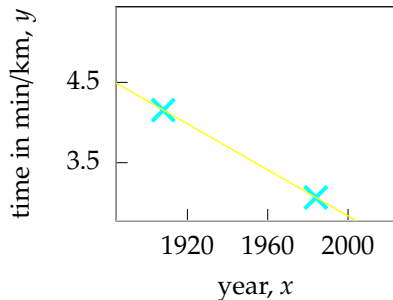
$$y_2 = mx_2 + c$$



Two Simultaneous Equations

A system of two simultaneous equations with two unknowns.

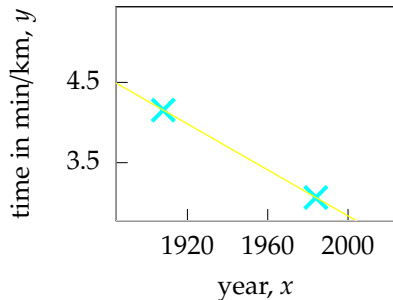
$$y_1 - y_2 = m(x_1 - x_2)$$



Two Simultaneous Equations

A system of two simultaneous equations with two unknowns.

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

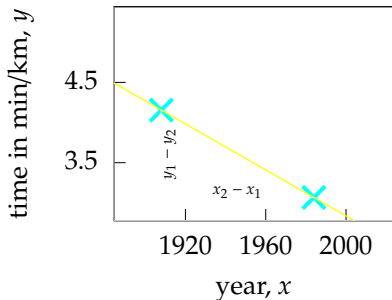


Two Simultaneous Equations

A system of two simultaneous equations with two unknowns.

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

$$c = y_1 - mx_1$$



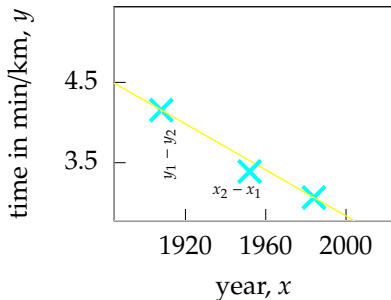
Two Simultaneous Equations

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

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

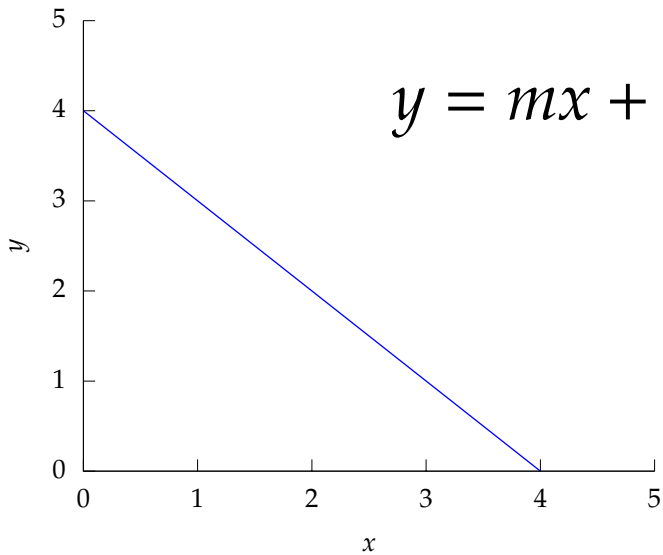
$$y_2 = mx_2 + c + \epsilon_2$$

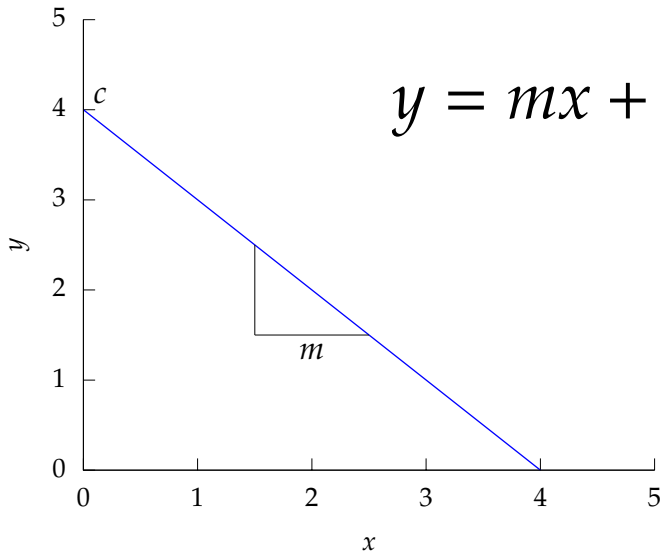
$$y_3 = mx_3 + c + \epsilon_3$$

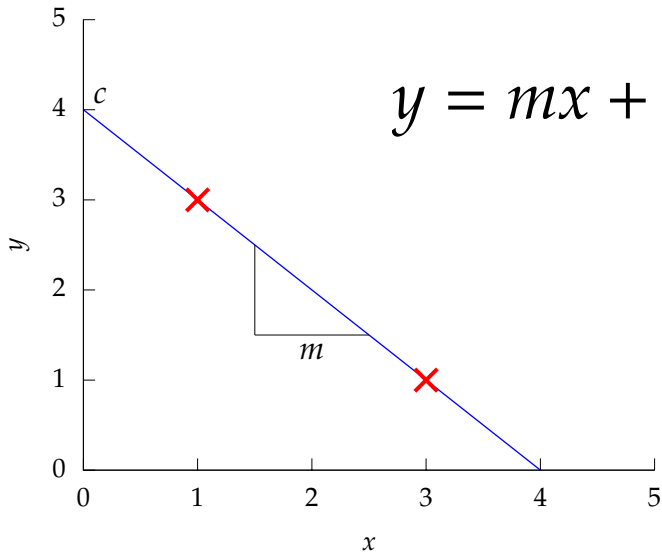
Noise Models

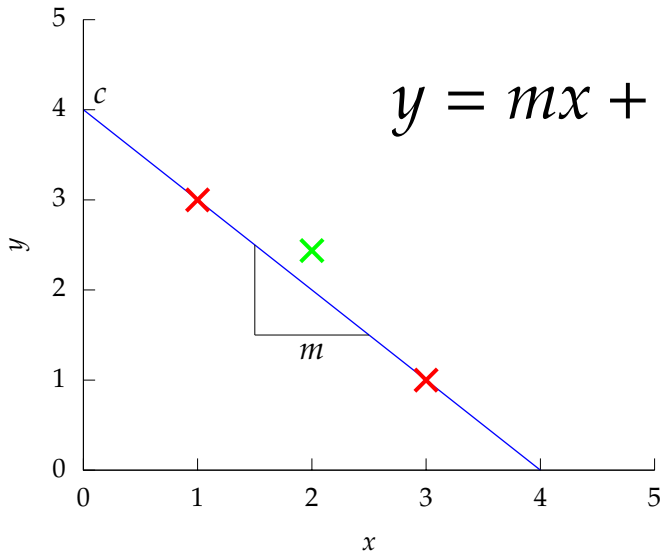
- ▶ We aren't modeling entire system.
- ▶ Noise model gives mismatch between model and data.
- ▶ Gaussian model justified by appeal to central limit theorem.
- ▶ Other models also possible (Student- t for heavy tails).
- ▶ Maximum likelihood with Gaussian noise leads to *least squares*.

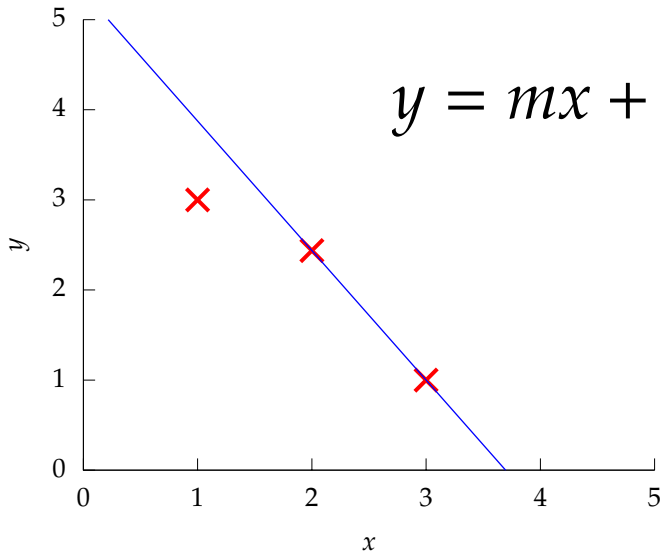
$$y = mx + c$$

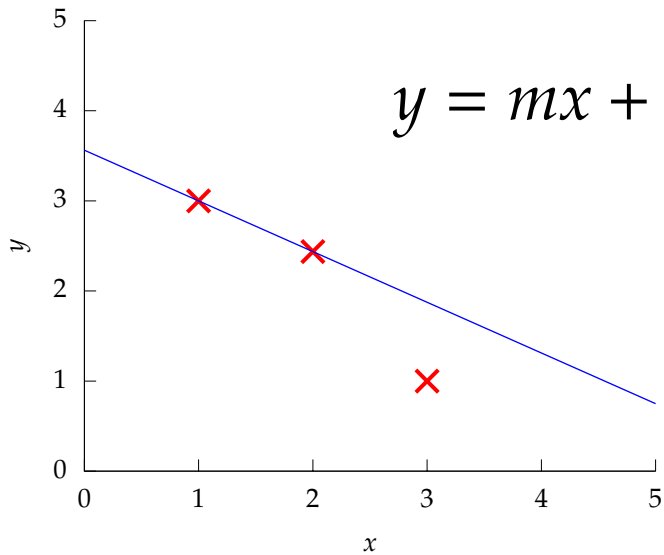


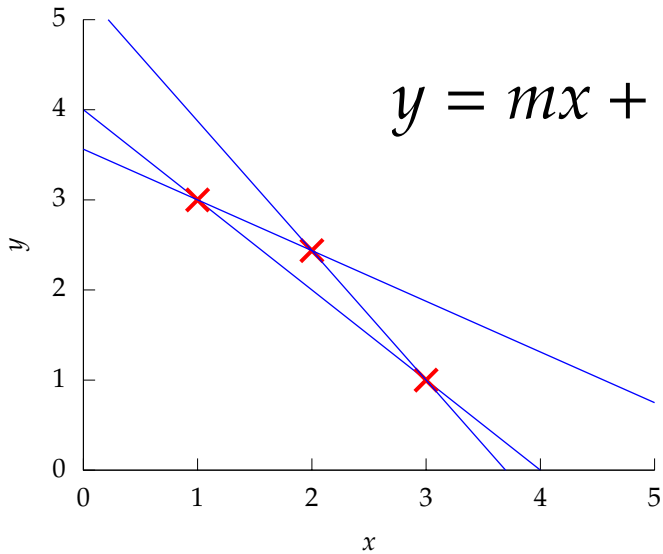












$$y = mx + c$$

point 1: $x = 1, y = 3$

$$3 = m + c$$

point 2: $x = 3, y = 1$

$$1 = 3m + c$$

point 3: $x = 2, y = 2.5$

$$2.5 = 2m + c$$



riens. L'opinion contraire est une illusion de l'esprit qui, perdant de vue les raisons fugitives du choix de la volonté dans les choses indifférentes, se persuade qu'elle s'est déterminée d'elle-même et sans motifs.

Nous devons donc envisager l'état présent de l'univers, comme l'effet de son état antérieur, et comme la cause de celui qui va suivre. Une intelligence qui, pour un instant donné, connaîtrait toutes les forces dont la nature est animée, et la situation respective des êtres qui la composent, si d'ailleurs elle était assez vaste pour soumettre ces données à l'analyse, embrasserait dans la même formule les mouvemens des plus grands corps de l'univers et ceux du plus léger atome : rien ne serait incertain pour elle, et l'avenir comme le passé, serait présent à ses yeux. L'esprit humain offre, dans la perfection qu'il a su donner à l'Astronomie, une faible esquisse de cette intelligence. Ses découvertes en Mécanique et en Géométrie, jointes à celle de la pesanteur universelle, l'ont mis à portée de comprendre dans les mêmes expressions analytiques, les états passés et futurs du système du monde. En appliquant la même méthode à quelques autres objets de ses connaissances, il est parvenu à ramener à des lois générales, les phénomènes observés, et à prévoir ceux que des circonstances données doivent faire éclore. Tous ces efforts dans la recherche de la vérité, tendent à le rapprocher sans cesse de l'intelligence que nous venons de concevoir, mais dont il restera toujours infiniment éloigné. Cette tendance propre à l'espèce humaine, est ce qui la rend supérieure aux animaux; et ses progrès en ce genre, distinguent les nations et les siècles, et font leur véritable gloire.

Rappelons-nous qu'autrefois, et à une époque qui

other, we say that its choice is an effect without a cause. It is then, says Leibnitz, the blind chance of the Epicureans. The contrary opinion is an illusion of the mind, which, losing sight of the evasive reasons of the choice of the will in indifferent things, believes that choice is determined of itself and without motives.

We ought then to regard the present state of the universe as the effect of its anterior state and as the cause of the one which is to follow. Given for one instant an intelligence which could comprehend all the forces by which nature is animated and the respective situation of the beings who compose it—an intelligence sufficiently vast to submit these data to analysis—it would embrace in the same formula the movements of the greatest bodies of the universe and those of the lightest atom; for it, nothing would be uncertain and the future, as the past, would be present to its eyes. The human mind offers, in the perfection which it has been able to give to astronomy, a feeble idea of this intelligence. Its discoveries in mechanics and geometry, added to that of universal gravity, have enabled it to comprehend in the same analytical expressions the past and future states of the system of the world. Applying the same method to some other objects of its knowledge, it has succeeded in referring to general laws observed phenomena and in foreseeing those which given circumstances ought to produce. All these efforts in the search for truth tend to lead it back continually to the vast intelligence which we have just mentioned, but from which it will always remain infinitely removed. This tendency, peculiar to the human race, is that which renders it superior to animals; and their progress

height: "The day will come when, by study pursued through several ages, the things now concealed will appear with evidence; and posterity will be astonished that truths so clear had escaped us." Clairaut then undertook to submit to analysis the perturbations which the comet had experienced by the action of the two great planets, Jupiter and Saturn; after immense calculations he fixed its next passage at the perihelion toward the beginning of April, 1759, which was actually verified by observation. The regularity which astronomy shows us in the movements of the comets doubtless exists also in all phenomena. .

The curve described by a simple molecule of air or vapor is regulated in a manner just as certain as the planetary orbits; the only difference between them is that which comes from our ignorance.

Probability is relative, in part to this ignorance, in part to our knowledge. We know that of three or a greater number of events a single one ought to occur; but nothing induces us to believe that one of them will occur rather than the others. In this state of indecision it is impossible for us to announce their occurrence with certainty. It is, however, probable that one of these events, chosen at will, will not occur because we see several cases equally possible which exclude its occurrence, while only a single one favors it.

The theory of chance consists in reducing all the events of the same kind to a certain number of cases equally possible, that is to say, to such as we may be equally undecided about in regard to their existence, and in determining the number of cases favorable to the event whose probability is sought. The ratio of

$$y = mx + c + \epsilon$$

point 1: $x = 1, y = 3$

$$3 = m + c + \epsilon_1$$

point 2: $x = 3, y = 1$

$$1 = 3m + c + \epsilon_2$$

point 3: $x = 2, y = 2.5$

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

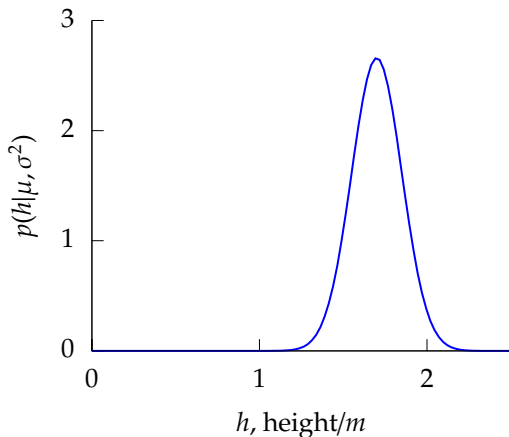
The Gaussian Density

- ▶ Perhaps the most common probability density.

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

- ▶ 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.

Two Important Gaussian Properties

Sum of Gaussians

- ▶ Sum of Gaussian variables is also Gaussian.

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

Two Important Gaussian Properties

Sum of Gaussians

- ▶ Sum of Gaussian variables is also Gaussian.

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

And the sum is distributed as

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

Two Important Gaussian Properties

Sum of Gaussians

- ▶ Sum of Gaussian variables is also Gaussian.

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

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].)

Two Important Gaussian Properties

Sum of Gaussians

- ▶ Sum of Gaussian variables is also Gaussian.

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

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].)

Two Important Gaussian Properties

Scaling a Gaussian

- ▶ Scaling a Gaussian leads to a Gaussian.

Two Important Gaussian Properties

Scaling a Gaussian

- ▶ Scaling a Gaussian leads to a Gaussian.

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

Two Important Gaussian Properties

Scaling a Gaussian

- ▶ Scaling a Gaussian leads to a Gaussian.

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

And the scaled density is distributed as

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

Outline

Quick Review: Overdetermined Systems

Underdetermined Systems

Bayesian Regression

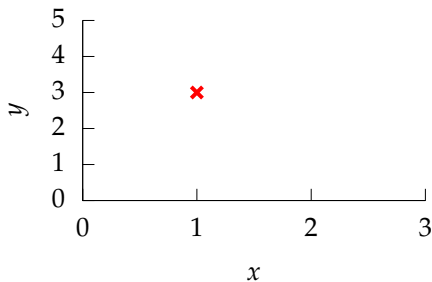
Univariate Bayesian Linear Regression

Bayesian Polynomials

Underdetermined System

What about two unknowns and *one* observation?

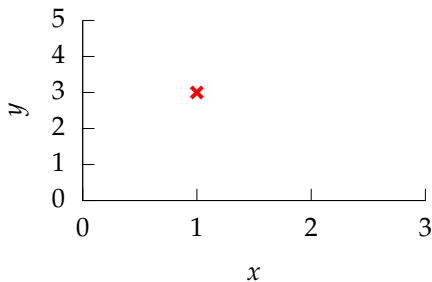
$$y_1 = mx_1 + c$$



Underdetermined System

Can compute m given c .

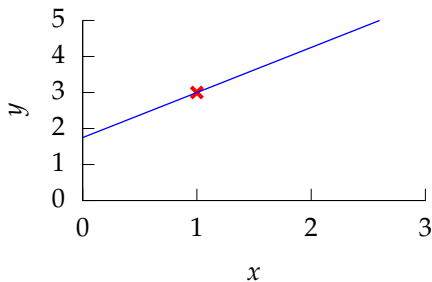
$$m = \frac{y_1 - c}{x}$$



Underdetermined System

Can compute m given c .

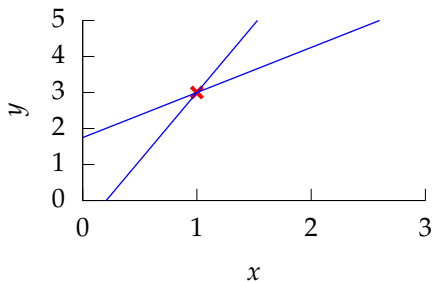
$$c = 1.75 \implies m = 1.25$$



Underdetermined System

Can compute m given c .

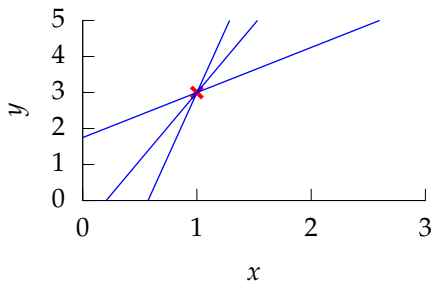
$$c = -0.777 \implies m = 3.78$$



Underdetermined System

Can compute m given c .

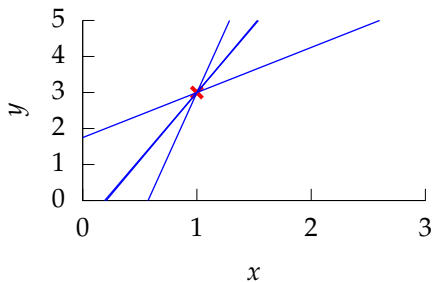
$$c = -4.01 \implies m = 7.01$$



Underdetermined System

Can compute m given c .

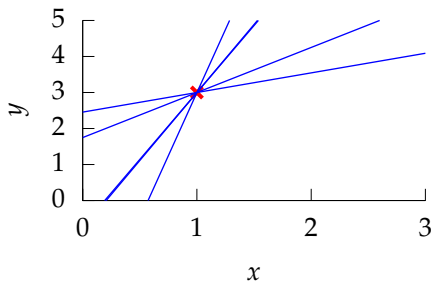
$$c = -0.718 \implies m = 3.72$$



Underdetermined System

Can compute m given c .

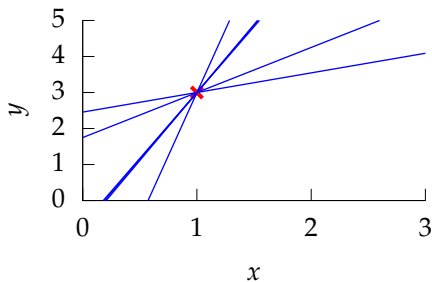
$$c = 2.45 \implies m = 0.545$$



Underdetermined System

Can compute m given c .

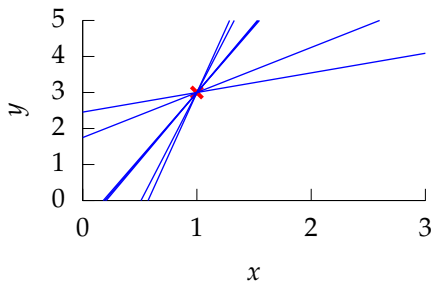
$$c = -0.657 \implies m = 3.66$$



Underdetermined System

Can compute m given c .

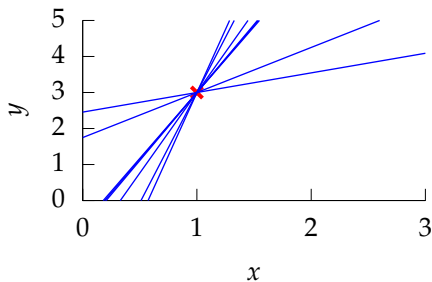
$$c = -3.13 \implies m = 6.13$$



Underdetermined System

Can compute m given c .

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



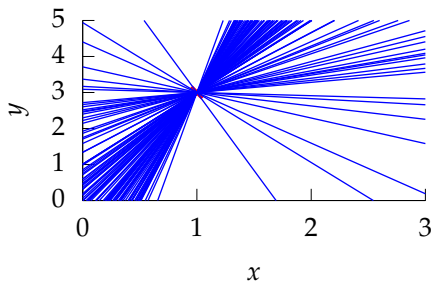
Underdetermined System

Can compute m given c .

Assume

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

we find a distribution of solutions.



Different Types of Uncertainty

- ▶ The first type of uncertainty we are assuming is *aleatoric* uncertainty.
- ▶ The second type of uncertainty we are assuming is *epistemic* uncertainty.

Aleatoric Uncertainty

- ▶ This is uncertainty we couldn't know even if we wanted to.
e.g. the result of a football match before it's played.
- ▶ Where a sheet of paper might land on the floor.

Outline

Quick Review: Overdetermined Systems

Underdetermined Systems

Bayesian Regression

Univariate Bayesian Linear Regression

Bayesian Polynomials

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

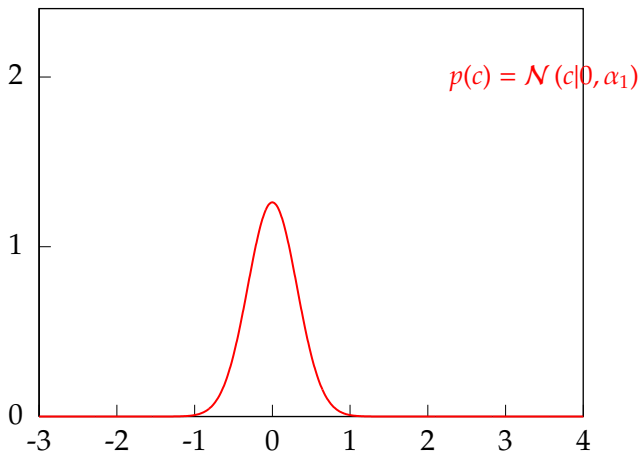


Figure : A Gaussian prior combined with a Gaussian likelihood for a Gaussian posterior.

Bayes Update

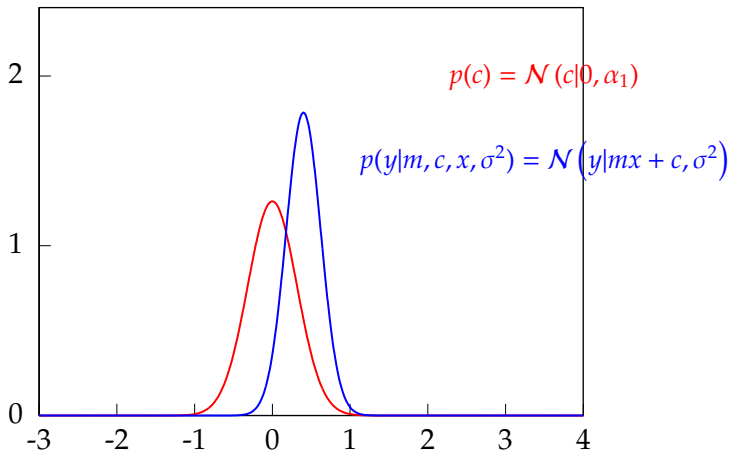


Figure : A Gaussian prior combined with a Gaussian likelihood for a Gaussian posterior.

Bayes Update

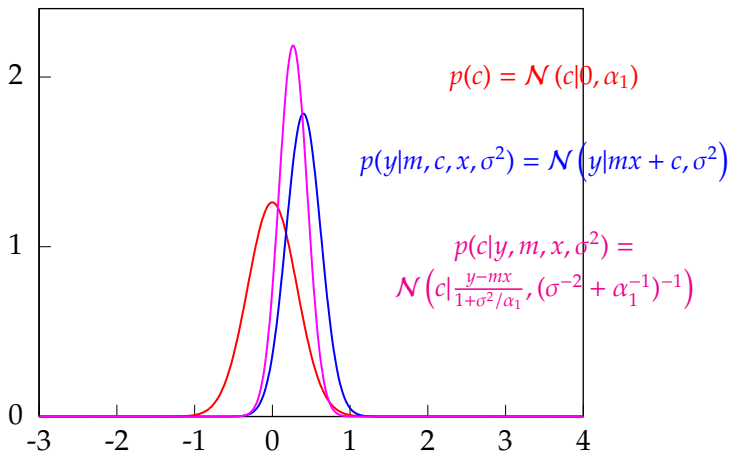


Figure : A Gaussian prior combined with a Gaussian likelihood for a Gaussian posterior.

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.

Main Trick

$$p(c) = \frac{1}{\sqrt{2\pi\alpha_1}} \exp\left(-\frac{1}{2\alpha_1}c^2\right)$$

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

Main Trick

$$p(c) = \frac{1}{\sqrt{2\pi\alpha_1}} \exp\left(-\frac{1}{2\alpha_1}c^2\right)$$

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

$$p(c|\mathbf{y}, \mathbf{x}, m, \sigma^2) = \frac{p(\mathbf{y}|\mathbf{x}, c, m, \sigma^2)p(c)}{p(\mathbf{y}|\mathbf{x}, m, \sigma^2)}$$

Main Trick

$$p(c) = \frac{1}{\sqrt{2\pi\alpha_1}} \exp\left(-\frac{1}{2\alpha_1}c^2\right)$$

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

$$p(c|\mathbf{y}, \mathbf{x}, m, \sigma^2) = \frac{p(\mathbf{y}|\mathbf{x}, c, m, \sigma^2)p(c)}{\int p(\mathbf{y}|\mathbf{x}, c, m, \sigma^2)p(c)dc}$$

Main Trick

$$p(c) = \frac{1}{\sqrt{2\pi\alpha_1}} \exp\left(-\frac{1}{2\alpha_1}c^2\right)$$

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

$$p(c|\mathbf{y}, \mathbf{x}, m, \sigma^2) \propto p(\mathbf{y}|\mathbf{x}, c, m, \sigma^2)p(c)$$

$$\begin{aligned}
\log p(c|\mathbf{y}, \mathbf{x}, m, \sigma^2) &= -\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - c - mx_i)^2 - \frac{1}{2\alpha_1} c^2 + \text{const} \\
&= -\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - mx_i)^2 - \left(\frac{n}{2\sigma^2} + \frac{1}{2\alpha_1} \right) c^2 \\
&\quad + c \frac{\sum_{i=1}^n (y_i - mx_i)}{\sigma^2},
\end{aligned}$$

complete the square of the quadratic form to obtain

$$\log p(c|\mathbf{y}, \mathbf{x}, m, \sigma^2) = -\frac{1}{2\tau^2} (c - \mu)^2 + \text{const},$$

where $\tau^2 = (n\sigma^{-2} + \alpha_1^{-1})^{-1}$ and $\mu = \frac{\tau^2}{\sigma^2} \sum_{n=1}^N (y_i - mx_i)$.

The Joint Density

- ▶ Really want to know the *joint* posterior density over the parameters c and m .
- ▶ Could now integrate out over m , but it's easier to consider the multivariate case.

Aleatoric Uncertainty

- ▶ This is uncertainty we couldn't know even if we wanted to.
e.g. the result of a football match before it's played.
- ▶ Where a sheet of paper might land on the floor.

Epistemic Uncertainty

- ▶ This is uncertainty we could in principle know the answer too. We just haven't observed enough yet, e.g. the result of a football match *after* it's played.
- ▶ What colour socks your lecturer is wearing.

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.

- ▶ Bayesian Inference

- ▶ 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.
- ▶ Bishop Section 1.2.3 (pg 21-24).
- ▶ Bishop Section 1.2.6 (start from just past eq 1.64 pg 30-32).

Outline

Quick Review: Overdetermined Systems

Underdetermined Systems

Bayesian Regression

Univariate Bayesian Linear Regression

Bayesian Polynomials

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

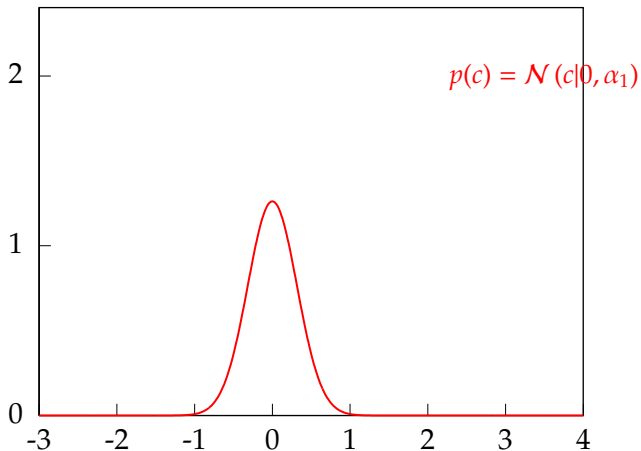


Figure : A Gaussian prior combined with a Gaussian likelihood for a Gaussian posterior.

Bayes Update

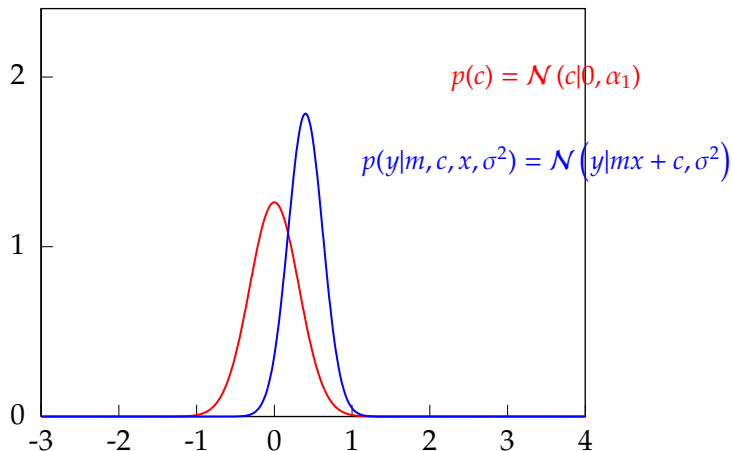


Figure : A Gaussian prior combined with a Gaussian likelihood for a Gaussian posterior.

Bayes Update

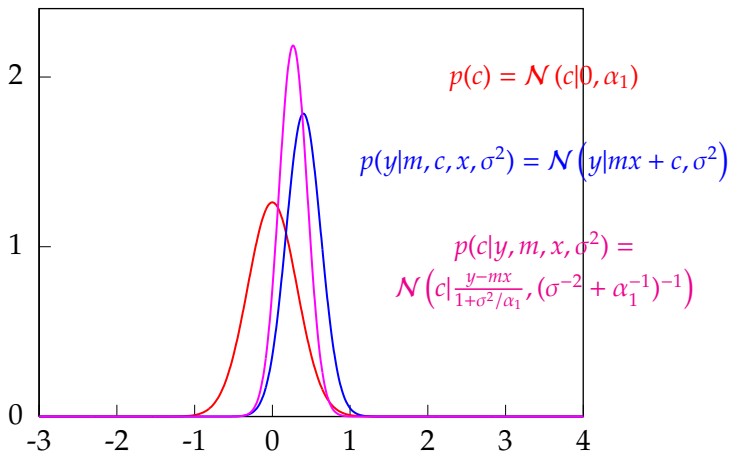


Figure : A Gaussian prior combined with a Gaussian likelihood for a Gaussian posterior.

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.

Main Trick

$$p(c) = \frac{1}{\sqrt{2\pi\alpha_1}} \exp\left(-\frac{1}{2\alpha_1}c^2\right)$$

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

Main Trick

$$p(c) = \frac{1}{\sqrt{2\pi\alpha_1}} \exp\left(-\frac{1}{2\alpha_1}c^2\right)$$

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

$$p(c|\mathbf{y}, \mathbf{x}, m, \sigma^2) = \frac{p(\mathbf{y}|\mathbf{x}, c, m, \sigma^2)p(c)}{p(\mathbf{y}|\mathbf{x}, m, \sigma^2)}$$

Main Trick

$$p(c) = \frac{1}{\sqrt{2\pi\alpha_1}} \exp\left(-\frac{1}{2\alpha_1}c^2\right)$$

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

$$p(c|\mathbf{y}, \mathbf{x}, m, \sigma^2) = \frac{p(\mathbf{y}|\mathbf{x}, c, m, \sigma^2)p(c)}{\int p(\mathbf{y}|\mathbf{x}, c, m, \sigma^2)p(c)dc}$$

Main Trick

$$p(c) = \frac{1}{\sqrt{2\pi\alpha_1}} \exp\left(-\frac{1}{2\alpha_1}c^2\right)$$

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

$$p(c|\mathbf{y}, \mathbf{x}, m, \sigma^2) \propto p(\mathbf{y}|\mathbf{x}, c, m, \sigma^2)p(c)$$

$$\begin{aligned}
\log p(c|\mathbf{y}, \mathbf{x}, m, \sigma^2) &= -\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - c - mx_i)^2 - \frac{1}{2\alpha_1} c^2 + \text{const} \\
&= -\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - mx_i)^2 - \left(\frac{n}{2\sigma^2} + \frac{1}{2\alpha_1} \right) c^2 \\
&\quad + c \frac{\sum_{i=1}^n (y_i - mx_i)}{\sigma^2},
\end{aligned}$$

complete the square of the quadratic form to obtain

$$\log p(c|\mathbf{y}, \mathbf{x}, m, \sigma^2) = -\frac{1}{2\tau^2} (c - \mu)^2 + \text{const},$$

where $\tau^2 = (n\sigma^{-2} + \alpha_1^{-1})^{-1}$ and $\mu = \frac{\tau^2}{\sigma^2} \sum_{n=1}^N (y_i - mx_i)$.

The Joint Density

- ▶ Really want to know the *joint* posterior density over the parameters c and m .
- ▶ Could now integrate out over m , but it's easier to consider the multivariate case.

Two Dimensional Gaussian

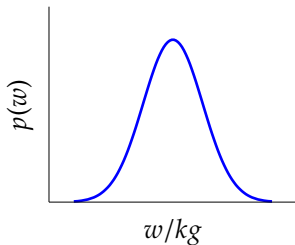
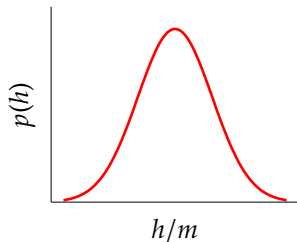
- ▶ Consider height, h/m and weight, w/kg .
- ▶ Could sample height from a distribution:

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

- ▶ And similarly weight:

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

Height and Weight Models

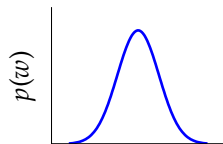
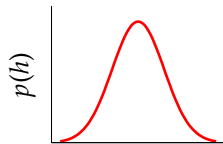
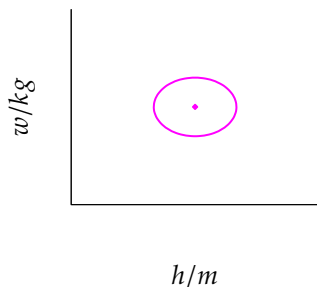


Gaussian distributions for height and weight.

Sampling Two Dimensional Variables

Marginal Distributions

Joint Distribution

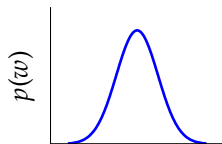
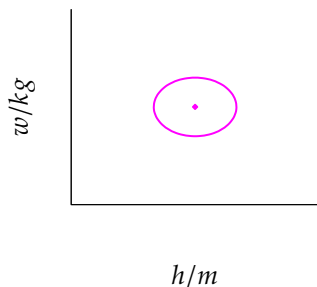


Samples of height and weight

Sampling Two Dimensional Variables

Marginal Distributions

Joint Distribution

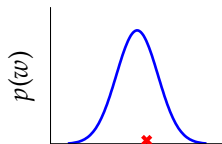
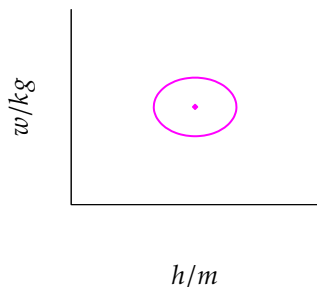


Samples of height and weight

Sampling Two Dimensional Variables

Marginal Distributions

Joint Distribution

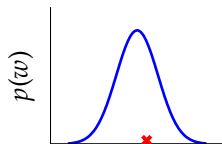
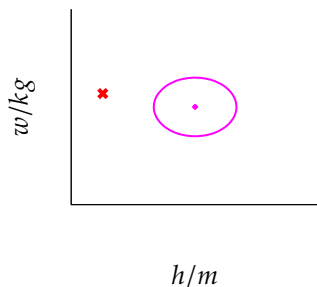


Samples of height and weight

Sampling Two Dimensional Variables

Marginal Distributions

Joint Distribution

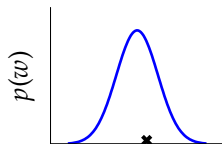
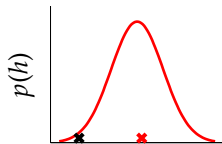
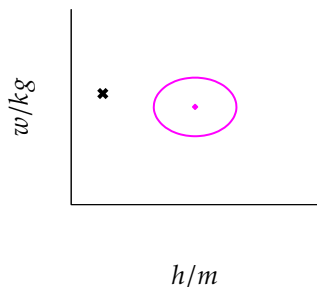


Samples of height and weight

Sampling Two Dimensional Variables

Marginal Distributions

Joint Distribution

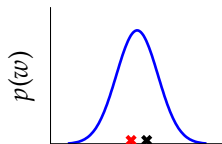
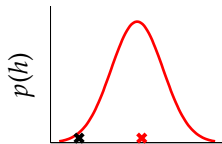
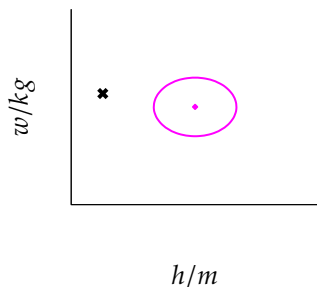


Samples of height and weight

Sampling Two Dimensional Variables

Marginal Distributions

Joint Distribution

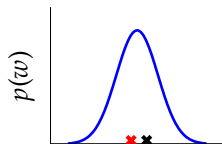
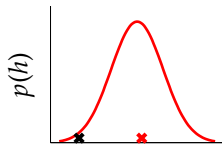
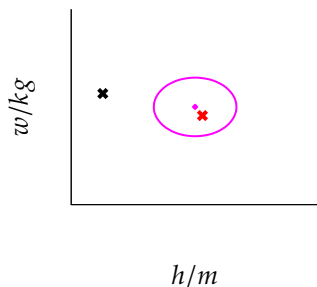


Samples of height and weight

Sampling Two Dimensional Variables

Marginal Distributions

Joint Distribution

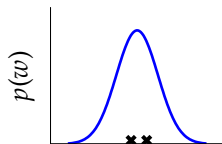
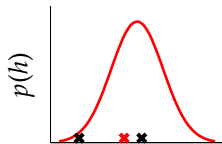
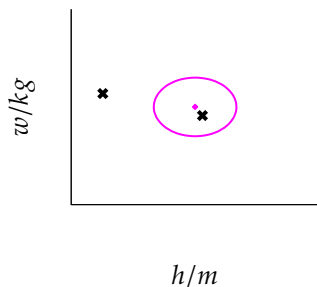


Samples of height and weight

Sampling Two Dimensional Variables

Marginal Distributions

Joint Distribution

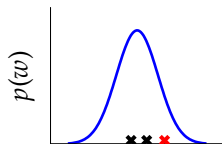
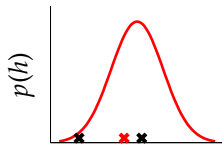
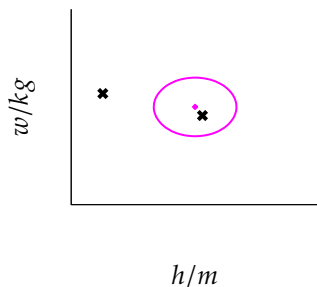


Samples of height and weight

Sampling Two Dimensional Variables

Marginal Distributions

Joint Distribution

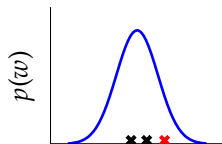
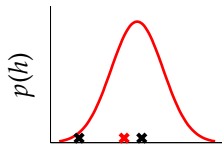
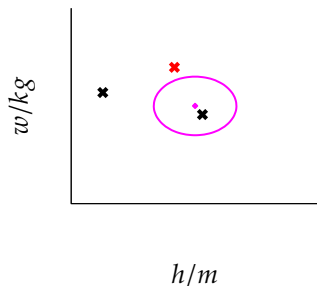


Samples of height and weight

Sampling Two Dimensional Variables

Marginal Distributions

Joint Distribution

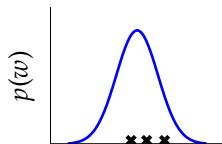
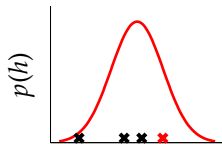
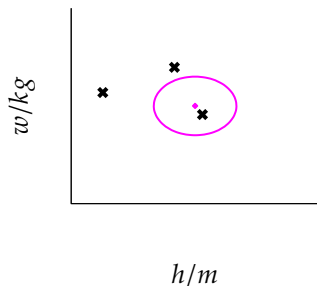


Samples of height and weight

Sampling Two Dimensional Variables

Marginal Distributions

Joint Distribution

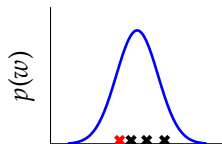
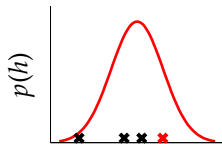
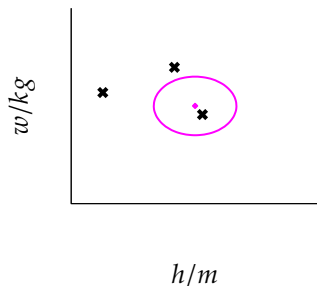


Samples of height and weight

Sampling Two Dimensional Variables

Marginal Distributions

Joint Distribution

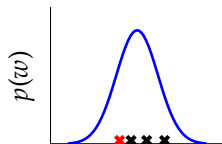
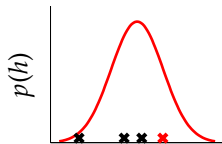
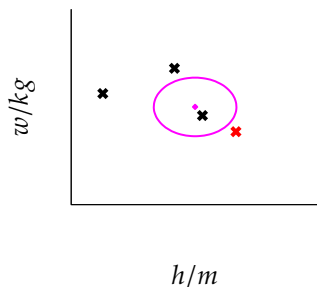


Samples of height and weight

Sampling Two Dimensional Variables

Marginal Distributions

Joint Distribution

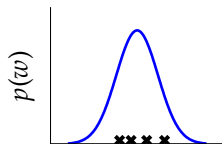
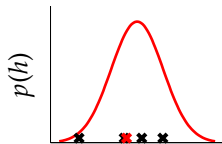
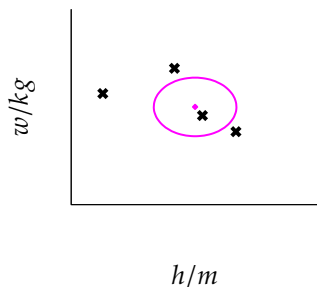


Samples of height and weight

Sampling Two Dimensional Variables

Marginal Distributions

Joint Distribution

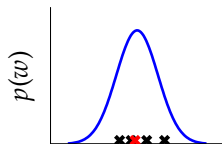
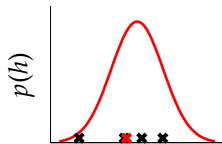
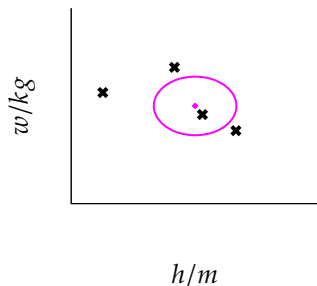


Samples of height and weight

Sampling Two Dimensional Variables

Marginal Distributions

Joint Distribution

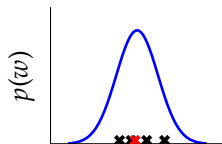
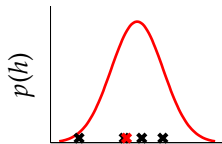
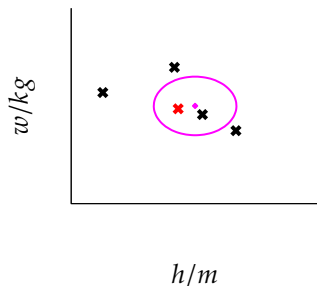


Samples of height and weight

Sampling Two Dimensional Variables

Marginal Distributions

Joint Distribution

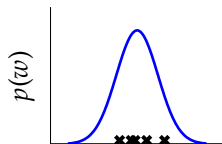
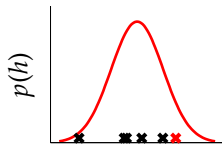
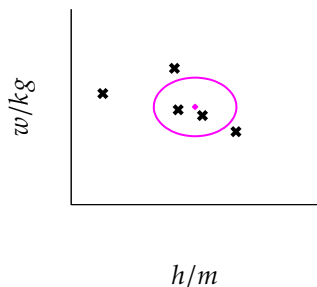


Samples of height and weight

Sampling Two Dimensional Variables

Marginal Distributions

Joint Distribution

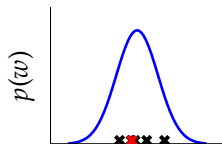
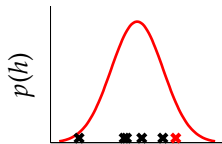
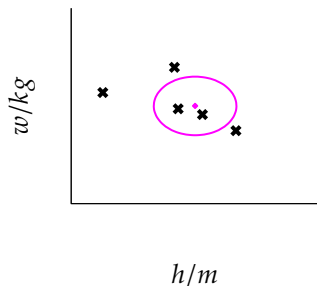


Samples of height and weight

Sampling Two Dimensional Variables

Marginal Distributions

Joint Distribution

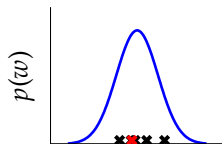
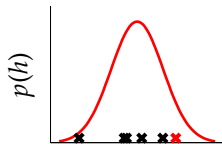
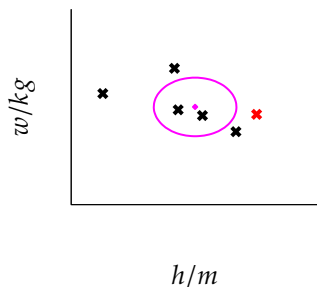


Samples of height and weight

Sampling Two Dimensional Variables

Marginal Distributions

Joint Distribution

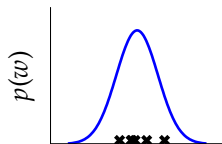
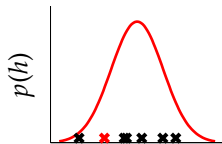
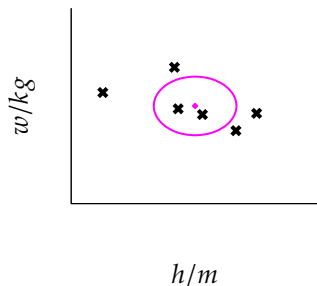


Samples of height and weight

Sampling Two Dimensional Variables

Marginal Distributions

Joint Distribution

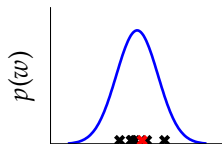
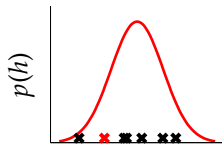
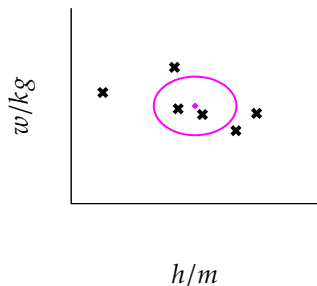


Samples of height and weight

Sampling Two Dimensional Variables

Marginal Distributions

Joint Distribution

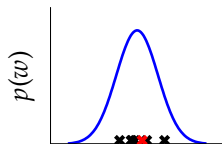
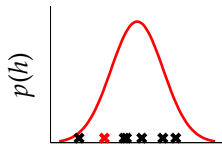
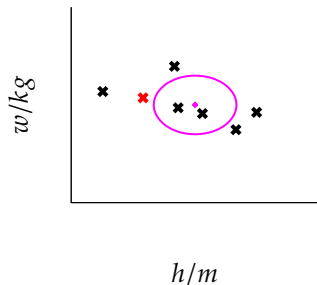


Samples of height and weight

Sampling Two Dimensional Variables

Marginal Distributions

Joint Distribution

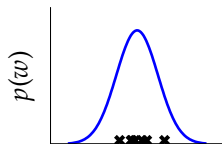
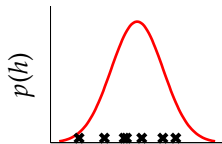
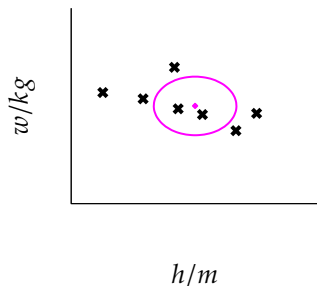


Samples of height and weight

Sampling Two Dimensional Variables

Marginal Distributions

Joint Distribution



Samples of height and weight

Independence Assumption

- ▶ This assumes height and weight are independent.

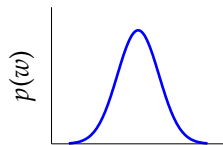
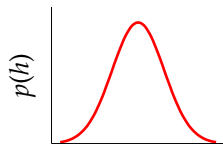
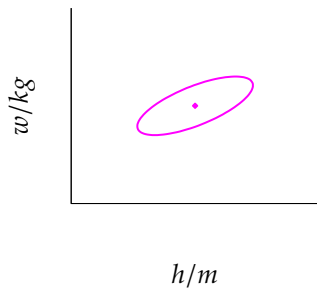
$$p(h, w) = p(h)p(w)$$

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

Sampling Two Dimensional Variables

Marginal Distributions

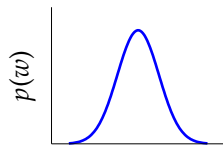
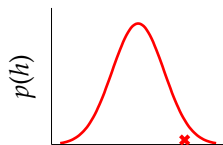
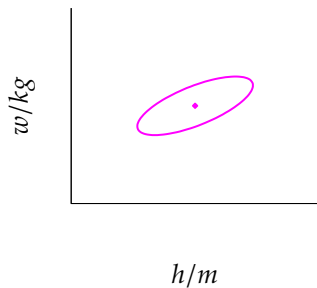
Joint Distribution



Sampling Two Dimensional Variables

Marginal Distributions

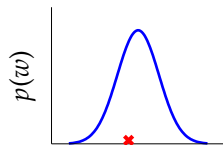
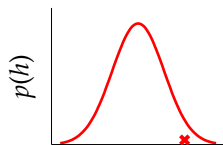
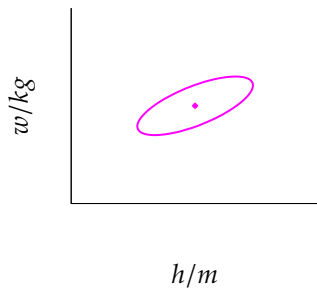
Joint Distribution



Sampling Two Dimensional Variables

Marginal Distributions

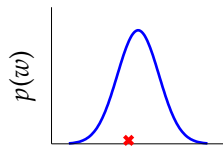
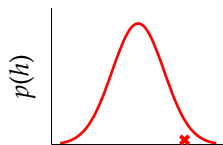
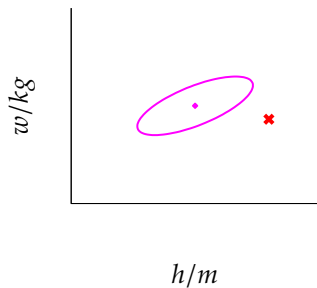
Joint Distribution



Sampling Two Dimensional Variables

Marginal Distributions

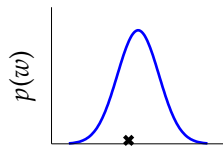
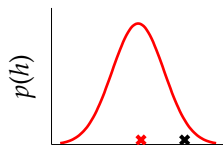
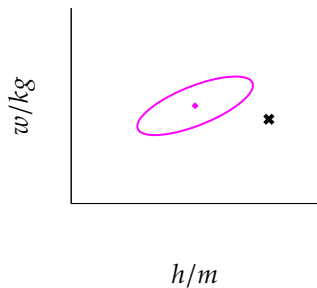
Joint Distribution



Sampling Two Dimensional Variables

Marginal Distributions

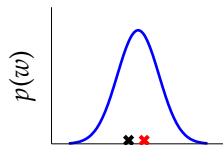
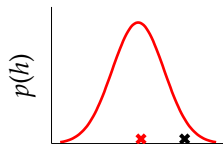
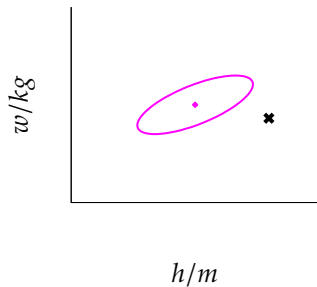
Joint Distribution



Sampling Two Dimensional Variables

Marginal Distributions

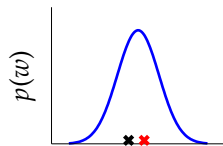
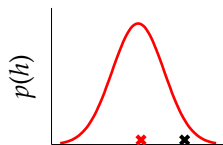
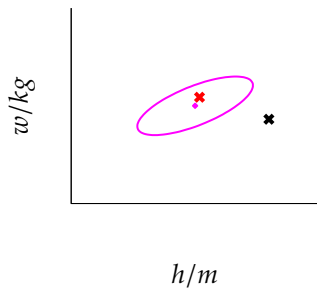
Joint Distribution



Sampling Two Dimensional Variables

Marginal Distributions

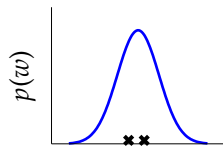
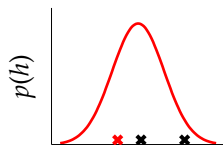
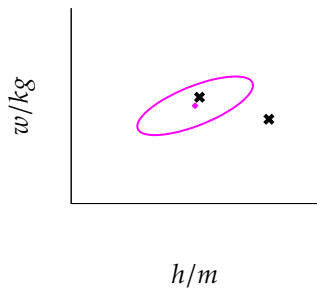
Joint Distribution



Sampling Two Dimensional Variables

Marginal Distributions

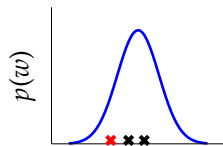
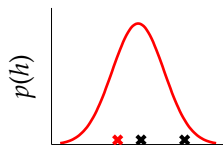
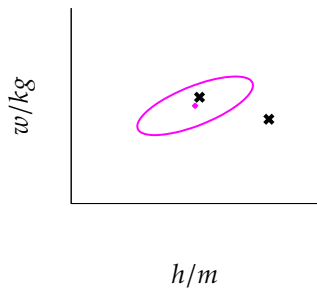
Joint Distribution



Sampling Two Dimensional Variables

Marginal Distributions

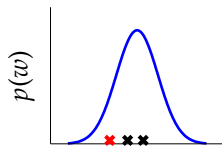
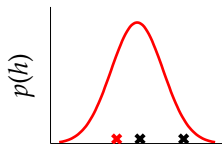
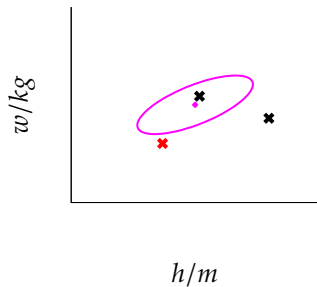
Joint Distribution



Sampling Two Dimensional Variables

Marginal Distributions

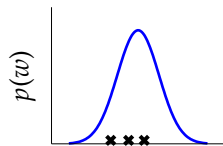
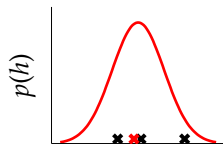
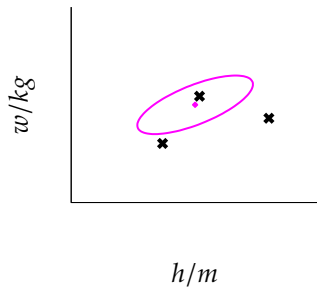
Joint Distribution



Sampling Two Dimensional Variables

Marginal Distributions

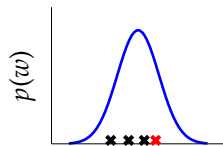
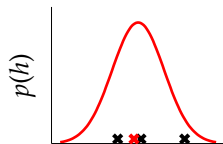
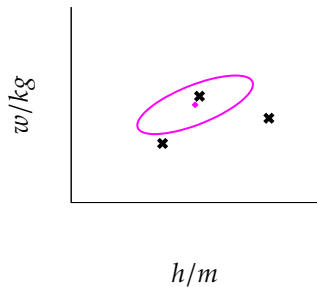
Joint Distribution



Sampling Two Dimensional Variables

Marginal Distributions

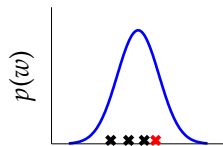
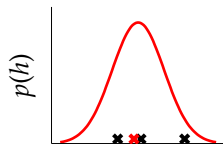
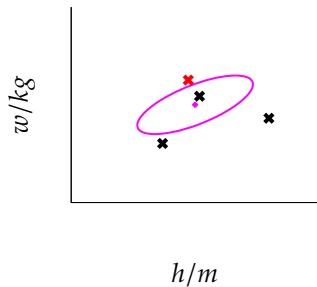
Joint Distribution



Sampling Two Dimensional Variables

Marginal Distributions

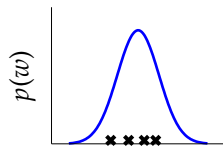
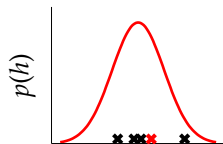
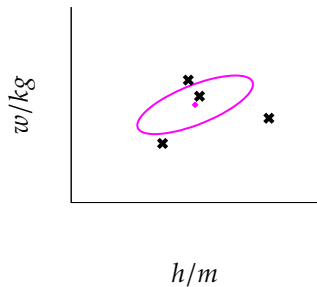
Joint Distribution



Sampling Two Dimensional Variables

Marginal Distributions

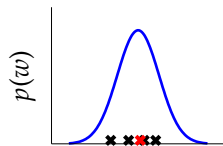
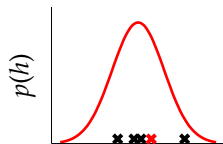
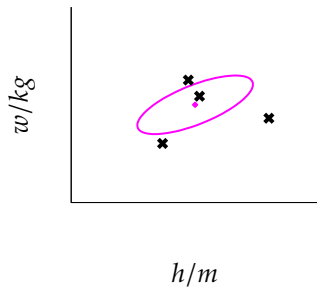
Joint Distribution



Sampling Two Dimensional Variables

Marginal Distributions

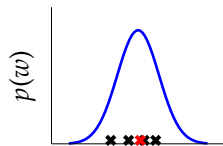
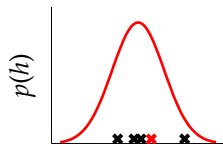
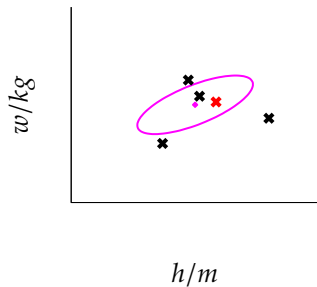
Joint Distribution



Sampling Two Dimensional Variables

Marginal Distributions

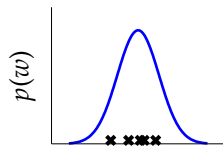
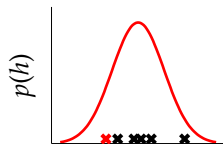
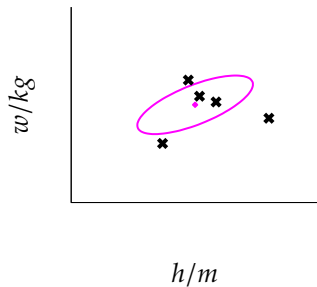
Joint Distribution



Sampling Two Dimensional Variables

Marginal Distributions

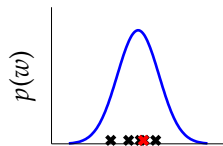
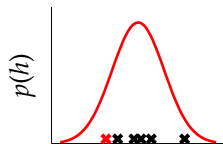
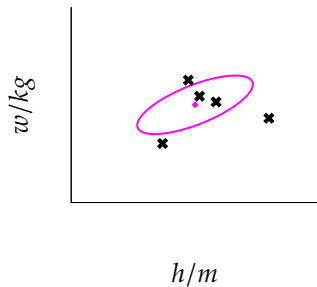
Joint Distribution



Sampling Two Dimensional Variables

Marginal Distributions

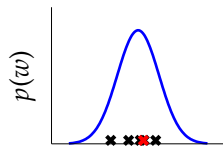
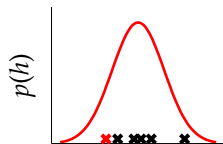
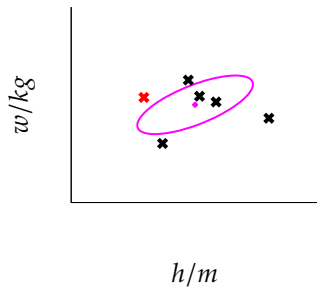
Joint Distribution



Sampling Two Dimensional Variables

Marginal Distributions

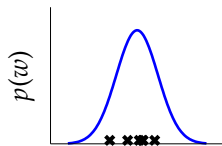
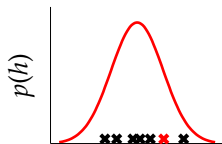
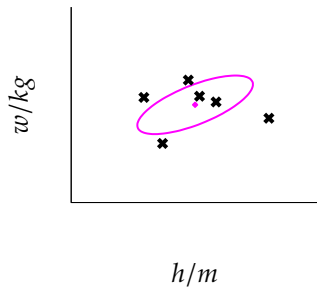
Joint Distribution



Sampling Two Dimensional Variables

Marginal Distributions

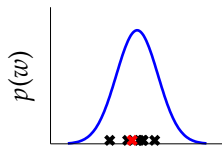
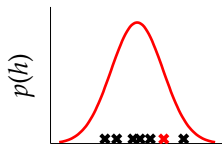
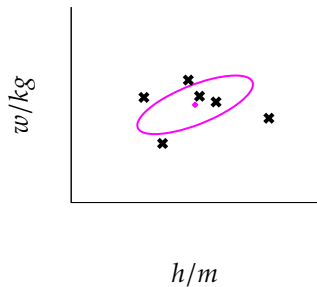
Joint Distribution



Sampling Two Dimensional Variables

Marginal Distributions

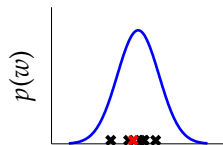
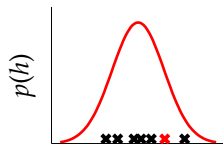
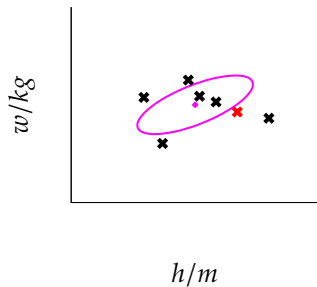
Joint Distribution



Sampling Two Dimensional Variables

Marginal Distributions

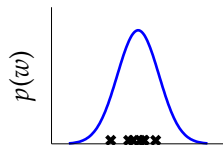
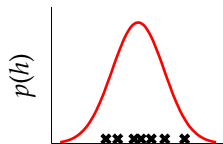
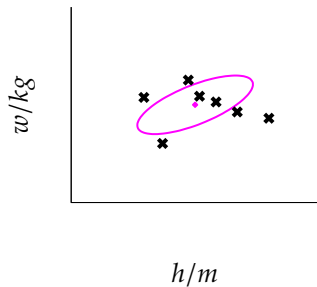
Joint Distribution



Sampling Two Dimensional Variables

Marginal Distributions

Joint Distribution



Independent Gaussians

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

Independent Gaussians

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

Independent Gaussians

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

Independent Gaussians

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

Correlated Gaussian

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

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

Correlated Gaussian

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

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

Correlated Gaussian

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:

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

Correlated Gaussian

Form correlated from original by rotating the data space using matrix \mathbf{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:

$$\mathbf{C} = \mathbf{R}\mathbf{D}\mathbf{R}^{\top}$$

Reading

- ▶ Section 2.3 of Bishop up to top of pg 85 (multivariate Gaussians).
- ▶ Section 3.3 of Bishop up to 159 (pg 152–159).

Outline

Quick Review: Overdetermined Systems

Underdetermined Systems

Bayesian Regression

Univariate Bayesian Linear Regression

Bayesian Polynomials

Revisit Olympics Data

- ▶ Use Bayesian approach on olympics data with polynomials.
- ▶ Choose a prior $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \alpha \mathbf{I})$ with $\alpha = 1$.
- ▶ Choose noise variance $\sigma^2 = 0.01$

Sampling the Prior

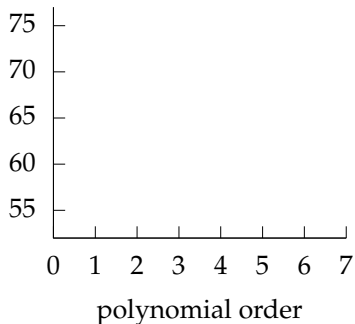
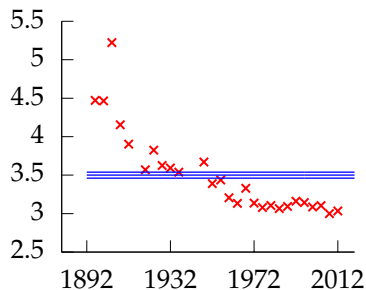
- ▶ Always useful to perform a ‘sanity check’ and sample from the prior before observing the data.
- ▶ Since $\mathbf{y} = \Phi\mathbf{w} + \epsilon$ just need to sample

$$w \sim \mathcal{N}(0, \alpha)$$

$$\epsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2)$$

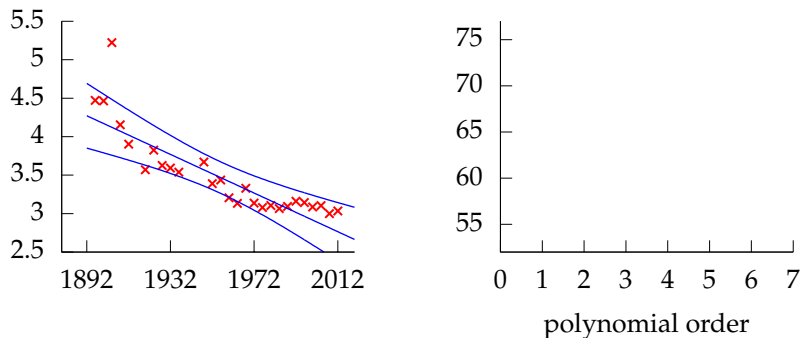
with $\alpha = 1$ and $\epsilon = 0.01$.

Polynomial Fits to Olympics Data



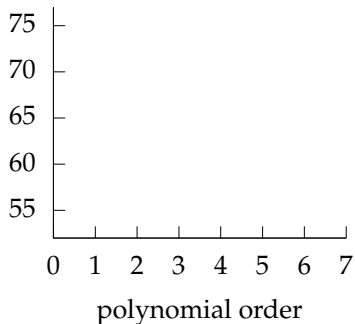
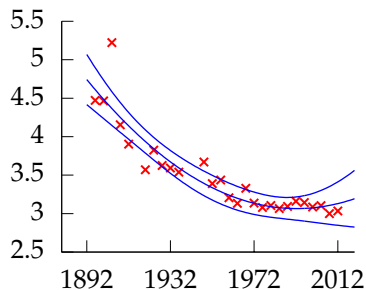
Left: fit to data, Right: marginal log likelihood. Polynomial order 0, model error 29.757, $\sigma^2 = 0.286$, $\sigma = 0.535$.

Polynomial Fits to Olympics Data



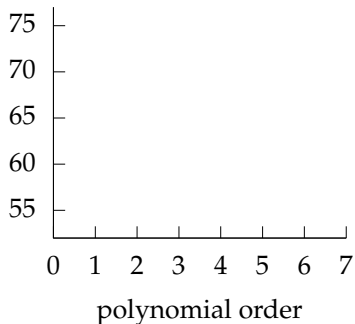
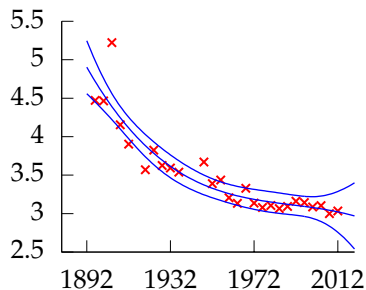
Left: fit to data, Right: marginal log likelihood. Polynomial order 1, model error 14.942, $\sigma^2 = 0.0749$, $\sigma = 0.274$.

Polynomial Fits to Olympics Data



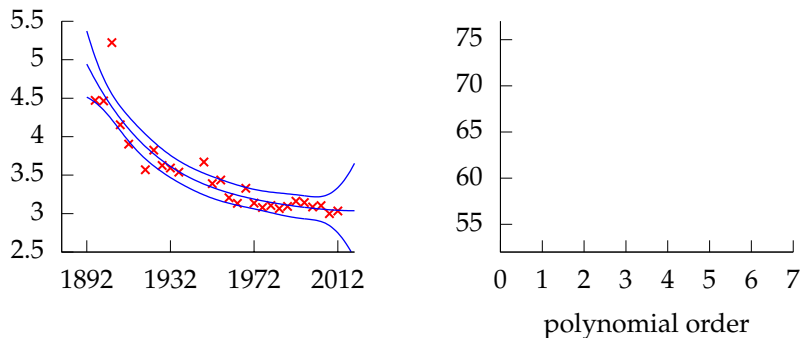
Left: fit to data, Right: marginal log likelihood. Polynomial order 2, model error 9.7206, $\sigma^2 = 0.0427$, $\sigma = 0.207$.

Polynomial Fits to Olympics Data



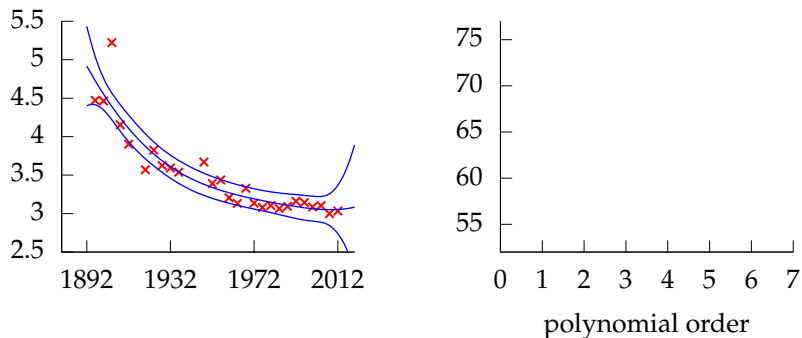
Left: fit to data, Right: marginal log likelihood. Polynomial order 3, model error 10.416, $\sigma^2 = 0.0402$, $\sigma = 0.200$.

Polynomial Fits to Olympics Data



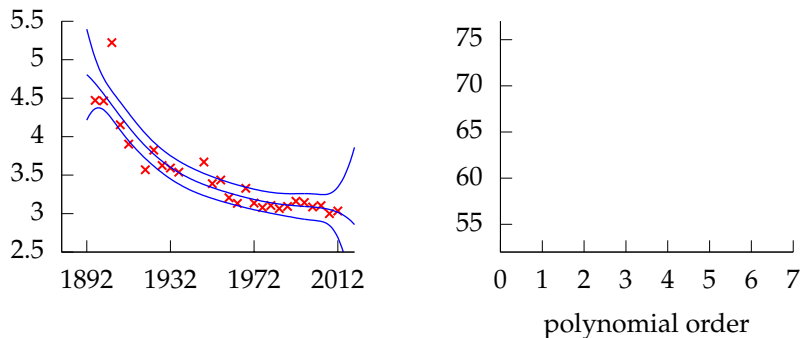
Left: fit to data, Right: marginal log likelihood. Polynomial order 4, model error 11.34, $\sigma^2 = 0.0401$, $\sigma = 0.200$.

Polynomial Fits to Olympics Data



Left: fit to data, Right: marginal log likelihood. Polynomial order 5, model error 11.986, $\sigma^2 = 0.0399$, $\sigma = 0.200$.

Polynomial Fits to Olympics Data

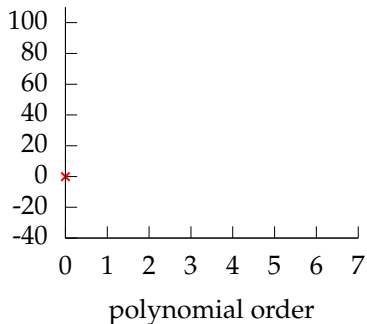
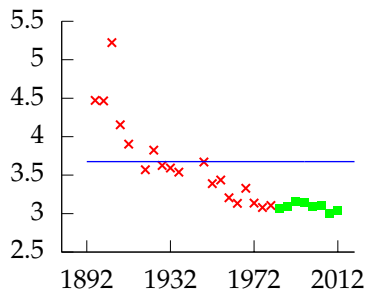


Left: fit to data, Right: marginal log likelihood. Polynomial order 6, model error 12.369, $\sigma^2 = 0.0384$, $\sigma = 0.196$.

Model Fit

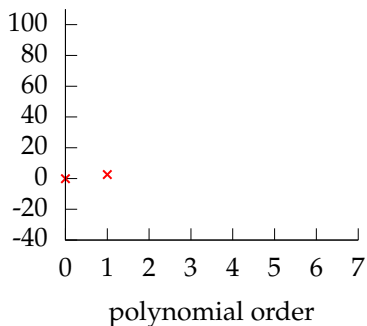
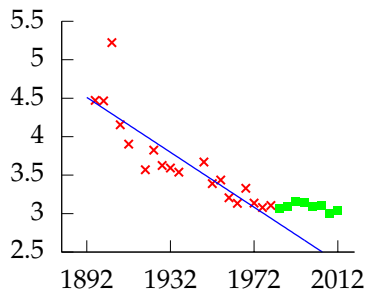
- ▶ Marginal likelihood doesn't always increase as model order increases.
- ▶ Bayesian model always has 2 parameters, regardless of how many basis functions (and here we didn't even fit them).
- ▶ Maximum likelihood model over fits through increasing number of parameters.
- ▶ Revisit maximum likelihood solution with validation set.

Recall: Validation Set for Maximum Likelihood



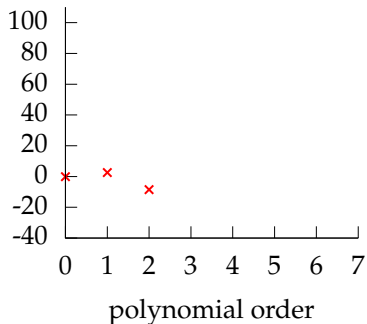
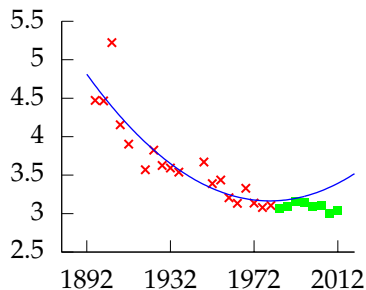
Left: fit to data, Right: model error. Polynomial order 0, training error -1.8774, validation error -0.13132, $\sigma^2 = 0.302$, $\sigma = 0.549$.

Recall: Validation Set for Maximum Likelihood



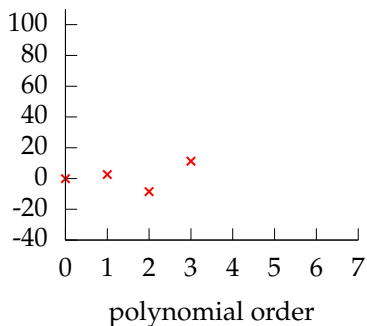
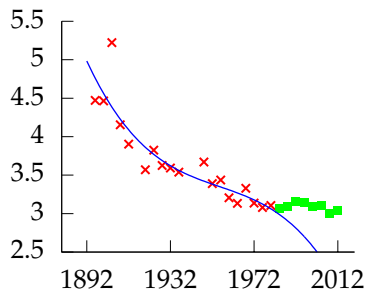
Left: fit to data, Right: model error. Polynomial order 1, training error -15.325, validation error 2.5863, $\sigma^2 = 0.0733$, $\sigma = 0.271$.

Recall: Validation Set for Maximum Likelihood



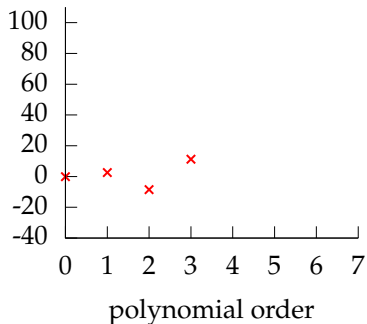
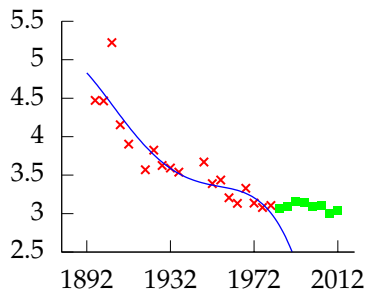
Left: fit to data, Right: model error. Polynomial order 2, training error -17.579, validation error -8.4831, $\sigma^2 = 0.0578$, $\sigma = 0.240$.

Recall: Validation Set for Maximum Likelihood



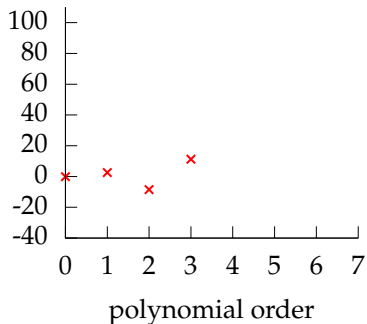
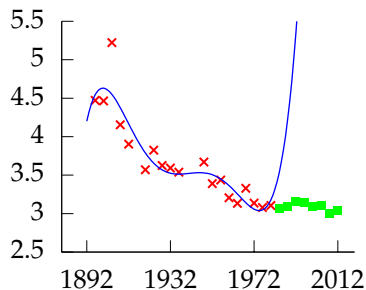
Left: fit to data, Right: model error. Polynomial order 3, training error -18.064, validation error 11.27, $\sigma^2 = 0.0549$, $\sigma = 0.234$.

Recall: Validation Set for Maximum Likelihood



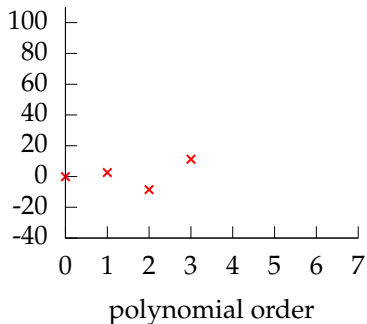
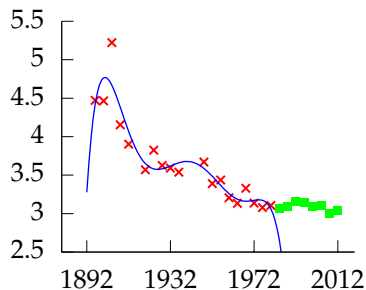
Left: fit to data, Right: model error. Polynomial order 4, training error -18.245, validation error 232.92, $\sigma^2 = 0.0539$, $\sigma = 0.232$.

Recall: Validation Set for Maximum Likelihood



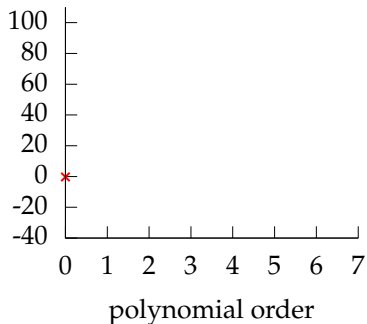
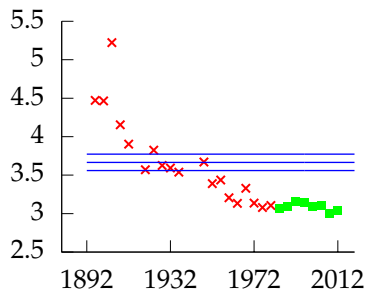
Left: fit to data, Right: model error. Polynomial order 5, training error -20.471, validation error 9898.1, $\sigma^2 = 0.0426$, $\sigma = 0.207$.

Recall: Validation Set for Maximum Likelihood



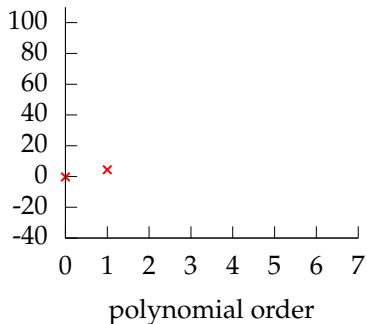
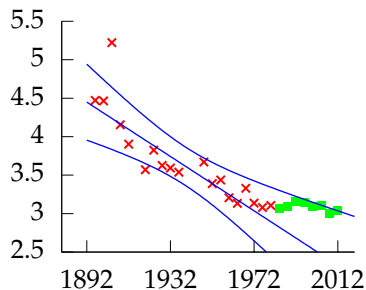
Left: fit to data, Right: model error. Polynomial order 6, training error -22.881, validation error 67775, $\sigma^2 = 0.0331$, $\sigma = 0.182$.

Validation Set



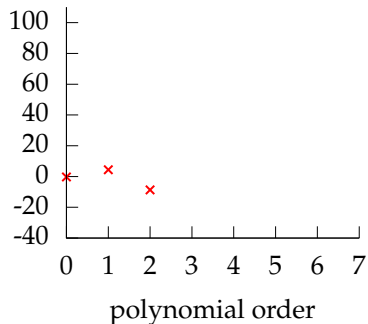
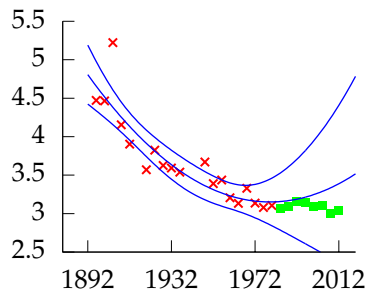
Left: fit to data, Right: model error. Polynomial order 0, training error 29.757, validation error -0.29243, $\sigma^2 = 0.302$, $\sigma = 0.550$.

Validation Set



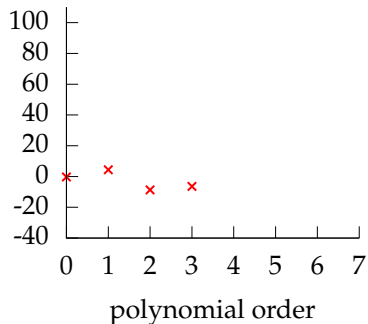
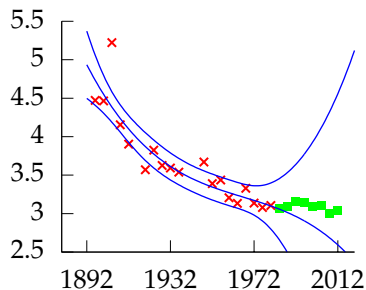
Left: fit to data, Right: model error. Polynomial order 1, training error 14.942, validation error 4.4027, $\sigma^2 = 0.0762$, $\sigma = 0.276$.

Validation Set



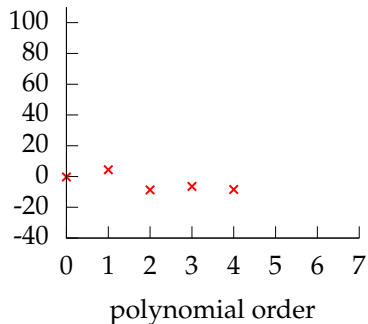
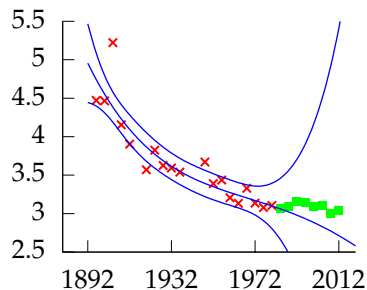
Left: fit to data, Right: model error. Polynomial order 2, training error 9.7206, validation error -8.6623, $\sigma^2 = 0.0580$, $\sigma = 0.241$.

Validation Set



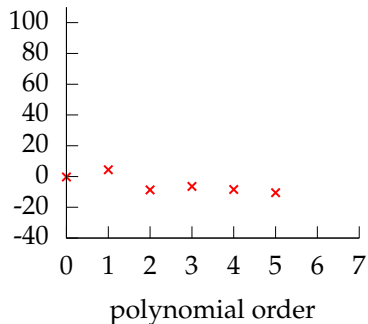
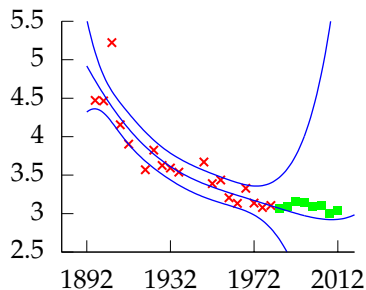
Left: fit to data, Right: model error. Polynomial order 3, training error 10.416, validation error -6.4726, $\sigma^2 = 0.0555$, $\sigma = 0.236$.

Validation Set



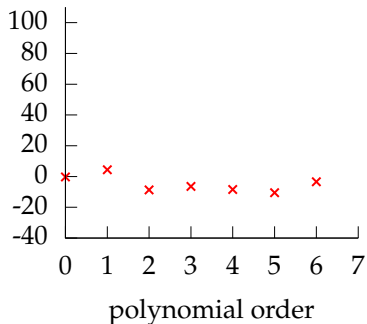
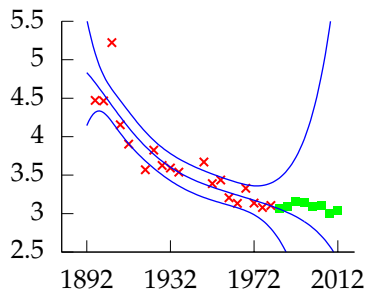
Left: fit to data, Right: model error. Polynomial order 4, training error 11.34, validation error -8.431, $\sigma^2 = 0.0555$, $\sigma = 0.236$.

Validation Set



Left: fit to data, Right: model error. Polynomial order 5, training error 11.986, validation error -10.483, $\sigma^2 = 0.0551$, $\sigma = 0.235$.

Validation Set



Left: fit to data, Right: model error. Polynomial order 6, training error 12.369, validation error -3.3823, $\sigma^2 = 0.0537$, $\sigma = 0.232$.

Regularized Mean

- ▶ Validation fit here based on mean solution for \mathbf{w} only.
- ▶ For Bayesian solution

$$\boldsymbol{\mu}_w = \left[\sigma^{-2} \boldsymbol{\Phi}^\top \boldsymbol{\Phi} + \alpha^{-1} \mathbf{I} \right]^{-1} \sigma^{-2} \boldsymbol{\Phi}^\top \mathbf{y}$$

instead of

$$\mathbf{w}^* = \left[\boldsymbol{\Phi}^\top \boldsymbol{\Phi} \right]^{-1} \boldsymbol{\Phi}^\top \mathbf{y}$$

- ▶ Two are equivalent when $\alpha \rightarrow \infty$.
- ▶ Equivalent to a prior for \mathbf{w} with infinite variance.
- ▶ In other cases $\alpha \mathbf{I}$ *regularizes* the system (keeps parameters smaller).

Sampling the Posterior

- ▶ Now check samples by extracting \mathbf{w} from the *posterior*.
- ▶ Now for $\mathbf{y} = \mathbf{\Phi}\mathbf{w} + \epsilon$ need

$$w \sim \mathcal{N}(\boldsymbol{\mu}_w, \mathbf{C}_w)$$

$$\text{with } \mathbf{C}_w = [\sigma^{-2}\mathbf{\Phi}^\top\mathbf{\Phi} + \alpha^{-1}\mathbf{I}]^{-1} \text{ and } \boldsymbol{\mu}_w = \mathbf{C}_w\sigma^{-2}\mathbf{\Phi}^\top\mathbf{y}$$

$$\epsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2)$$

with $\alpha = 1$ and $\epsilon = 0.01$.

Marginal Likelihood

- ▶ The marginal likelihood can also be computed, it has the form:

$$p(\mathbf{y}|\mathbf{X}, \sigma^2, \alpha) = \frac{1}{(2\pi)^{\frac{n}{2}} |\mathbf{K}|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}\mathbf{y}^\top \mathbf{K}^{-1} \mathbf{y}\right)$$

where $\mathbf{K} = \alpha \mathbf{\Phi} \mathbf{\Phi}^\top + \sigma^2 \mathbf{I}$.

- ▶ So it is a zero mean n -dimensional Gaussian with covariance matrix \mathbf{K} .

Computing the Expected Output

- ▶ Given the posterior for the parameters, how can we compute the expected output at a given location?
- ▶ Output of model at location \mathbf{x}_i is given by

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

- ▶ We want the expected output under the posterior density, $p(\mathbf{w}|\mathbf{y}, \mathbf{X}, \sigma^2, \alpha)$.
- ▶ Mean of mapping function will be given by

$$\begin{aligned}\langle f(\mathbf{x}_i; \mathbf{w}) \rangle_{p(\mathbf{w}|\mathbf{y}, \mathbf{X}, \sigma^2, \alpha)} &= \phi_i^\top \langle \mathbf{w} \rangle_{p(\mathbf{w}|\mathbf{y}, \mathbf{X}, \sigma^2, \alpha)} \\ &= \phi_i^\top \boldsymbol{\mu}_w\end{aligned}$$

Variance of Expected Output

- ▶ Variance of model at location \mathbf{x}_i is given by

$$\begin{aligned}\text{var}(f(\mathbf{x}_i; \mathbf{w})) &= \langle (f(\mathbf{x}_i; \mathbf{w}))^2 \rangle - \langle f(\mathbf{x}_i; \mathbf{w}) \rangle^2 \\ &= \boldsymbol{\phi}_i^\top \langle \mathbf{w}\mathbf{w}^\top \rangle \boldsymbol{\phi}_i - \boldsymbol{\phi}_i^\top \langle \mathbf{w} \rangle \langle \mathbf{w} \rangle^\top \boldsymbol{\phi}_i \\ &= \boldsymbol{\phi}_i^\top \mathbf{C}_i \boldsymbol{\phi}_i\end{aligned}$$

where all these expectations are taken under the posterior density, $p(\mathbf{w}|\mathbf{y}, \mathbf{X}, \sigma^2, \alpha)$.

Reading

- ▶ Section 3.7–3.8 of Rogers and Girolami (pg 122–133).
- ▶ Section 3.4 of Bishop (pg 161–165).

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 évènements. In *Mémoires de mathématique et de physique, présentés à l'Académie Royale des Sciences, par divers savans, & lû dans ses assemblées* 6, pages 621–656, 1774. Translated in Stigler (1986).
- P. S. Laplace. *Essai philosophique sur les probabilités*. Courcier, Paris, 2nd edition, 1814. Sixth edition of 1840 translated and reprinted (1951) as *A Philosophical Essay on Probabilities*, New York: Dover; fifth edition of 1825 reprinted 1986 with notes by Bernard Bru, Paris: Christian Bourgeois Éditeur, translated by Andrew Dale (1995) as *Philosophical Essay on Probabilities*, New York:Springer-Verlag.
- 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.