

# Linear regression

Sep 2022

Murphy chap7

## Linear regression - MSE

Typically we have a set of training data  $(x_1, y_1) \dots (x_N, y_N)$  from which to estimate the parameters  $\beta$ . Each  $x_i = (x_{i1}, x_{i2}, \dots, x_{ip})^T$  is a vector of feature measurements for the  $i$ th case. The most popular estimation method is *least squares*, in which we pick the coefficients  $\beta = (\beta_0, \beta_1, \dots, \beta_p)^T$  to minimize the residual sum of squares

$$\begin{aligned} \text{RSS}(\beta) &= \sum_{i=1}^N (y_i - f(x_i))^2 \\ &= \sum_{i=1}^N \left( y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2. \end{aligned} \tag{3.2}$$

How do we minimize (3.2)? Denote by  $\mathbf{X}$  the  $N \times (p + 1)$  matrix with each row an input vector (with a 1 in the first position), and similarly let  $\mathbf{y}$  be the  $N$ -vector of outputs in the training set. Then we can write the residual sum-of-squares as

$$\text{RSS}(\beta) = (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta). \tag{3.3}$$

## Linear regression - MSE

$$\text{RSS}(\beta) = (\mathbf{y} - \mathbf{X}\beta)^T(\mathbf{y} - \mathbf{X}\beta). \quad (3.3)$$

This is a quadratic function in the  $p + 1$  parameters. Differentiating with respect to  $\beta$  we obtain

$$\begin{aligned} \frac{\partial \text{RSS}}{\partial \beta} &= -2\mathbf{X}^T(\mathbf{y} - \mathbf{X}\beta) \\ \frac{\partial^2 \text{RSS}}{\partial \beta \partial \beta^T} &= 2\mathbf{X}^T \mathbf{X}. \end{aligned} \quad (3.4)$$

Assuming (for the moment) that  $\mathbf{X}$  has full column rank, and hence  $\mathbf{X}^T \mathbf{X}$  is positive definite, we set the first derivative to zero

$$\mathbf{X}^T(\mathbf{y} - \mathbf{X}\beta) = 0 \quad (3.5)$$

to obtain the unique solution

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}. \quad (3.6)$$

## Linear regression - MSE

We denote the column vectors of  $\mathbf{X}$  by  $\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_p$ , with  $\mathbf{x}_0 \equiv 1$ .

These vectors span a subspace of  $\mathbb{R}^N$ , also referred to as the column space of  $\mathbf{X}$ .

We minimize  $\text{RSS}(\beta) = \|\mathbf{y} - \mathbf{X}\beta\|^2$  by choosing  $\hat{\beta}$  so that the residual vector  $\mathbf{y} - \hat{\mathbf{y}}$  is orthogonal to this subspace.

$$\mathbf{X}^T(\mathbf{y} - \mathbf{X}\beta) = 0$$

# Linear regression - Geometry

We seek a vector  $\hat{\mathbf{y}} \in \mathbb{R}^N$  that lies in this linear subspace and is as close as possible to  $\mathbf{y}$ , i.e., we want to find

$$\operatorname{argmin}_{\hat{\mathbf{y}} \in \text{span}(\{\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_D\})} \|\mathbf{y} - \hat{\mathbf{y}}\|_2. \quad (7.18)$$

Since  $\hat{\mathbf{y}} \in \text{span}(\mathbf{X})$ , there exists some weight vector  $\mathbf{w}$  such that

$$\hat{\mathbf{y}} = w_1 \tilde{\mathbf{x}}_1 + \dots + w_D \tilde{\mathbf{x}}_D = \mathbf{X}\mathbf{w} \quad (7.19)$$

To minimize the norm of the residual,  $\mathbf{y} - \hat{\mathbf{y}}$ , we want the residual vector to be orthogonal to every column of  $\mathbf{X}$ , so  $\tilde{\mathbf{x}}_j^T (\mathbf{y} - \hat{\mathbf{y}}) = 0$  for  $j = 1 : D$ . Hence

$$\tilde{\mathbf{x}}_j^T (\mathbf{y} - \hat{\mathbf{y}}) = 0 \Rightarrow \mathbf{X}^T (\mathbf{y} - \mathbf{X}\mathbf{w}) = \mathbf{0} \Rightarrow \mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad (7.20)$$

Hence our projected value of  $\mathbf{y}$  is given by

$$\hat{\mathbf{y}} = \mathbf{X}\hat{\mathbf{w}} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad (7.21)$$

This corresponds to an **orthogonal projection** of  $\mathbf{y}$  onto the column space of  $\mathbf{X}$ . The projection matrix  $\mathbf{P} \triangleq \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$  is called the **hat matrix**, since it “puts the hat on  $\mathbf{y}$ ”.

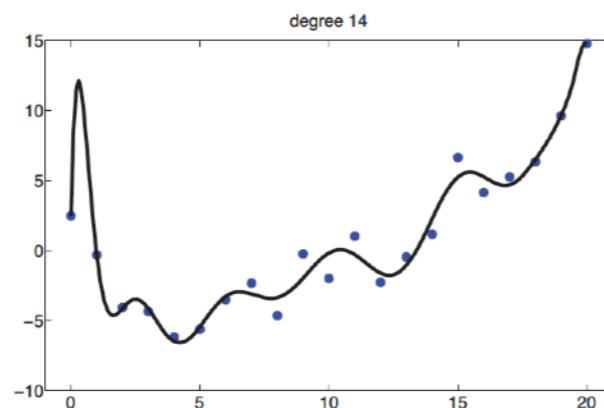
# Linear regression - Probability

## Linear regression

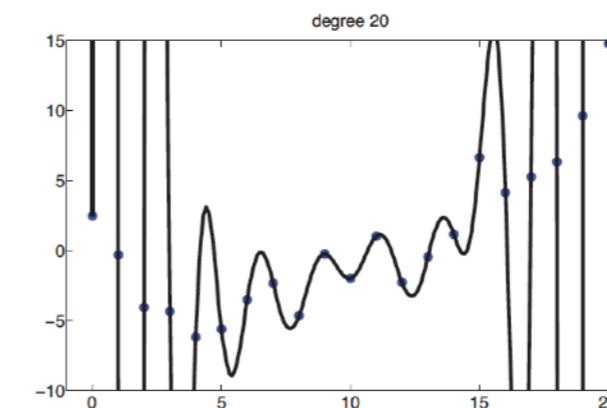
One of the most widely used models for regression is known as **linear regression**. This asserts that the response is a linear function of the inputs. This can be written as follows:

$$y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + \epsilon = \sum_{j=1}^D w_j x_j + \epsilon \quad (1.4)$$

where  $\mathbf{w}^T \mathbf{x}$  represents the inner or **scalar product** between the input vector  $\mathbf{x}$  and the model's **weight vector**  $\mathbf{w}$ <sup>7</sup>, and  $\epsilon$  is the **residual error** between our linear predictions and the true response.



(a)



(b)

**Figure 1.18** Polynomial of degrees 14 and 20 fit by least squares to 21 data points. Figure generated by linregPolyVsDegree.

# Linear regression - Probability

To make the connection between linear regression and Gaussians more explicit, we can rewrite the model in the following form:

$$p(y|\mathbf{x}, \boldsymbol{\theta}) = \mathcal{N}(y|\mu(\mathbf{x}), \sigma^2(\mathbf{x})) \quad (1.5)$$

This makes it clear that the model is a conditional probability density. In the simplest case, we assume  $\mu$  is a linear function of  $\mathbf{x}$ , so  $\mu = \mathbf{w}^T \mathbf{x}$ , and that the noise is fixed,  $\sigma^2(x) = \sigma^2$ . In this case,  $\boldsymbol{\theta} = (\mathbf{w}, \sigma^2)$  are the parameters of the model.

For example, suppose the input is 1 dimensional. We can represent the expected response as follows:

$$\mu(\mathbf{x}) = w_0 + w_1 x = \mathbf{w}^T \mathbf{x} \quad (1.6)$$

where  $w_0$  is the intercept or **bias** term,  $w_1$  is the slope, and where we have defined the vector  $\mathbf{x} = (1, x)$ . (Prepending a constant 1 term to an input vector is a common notational trick which

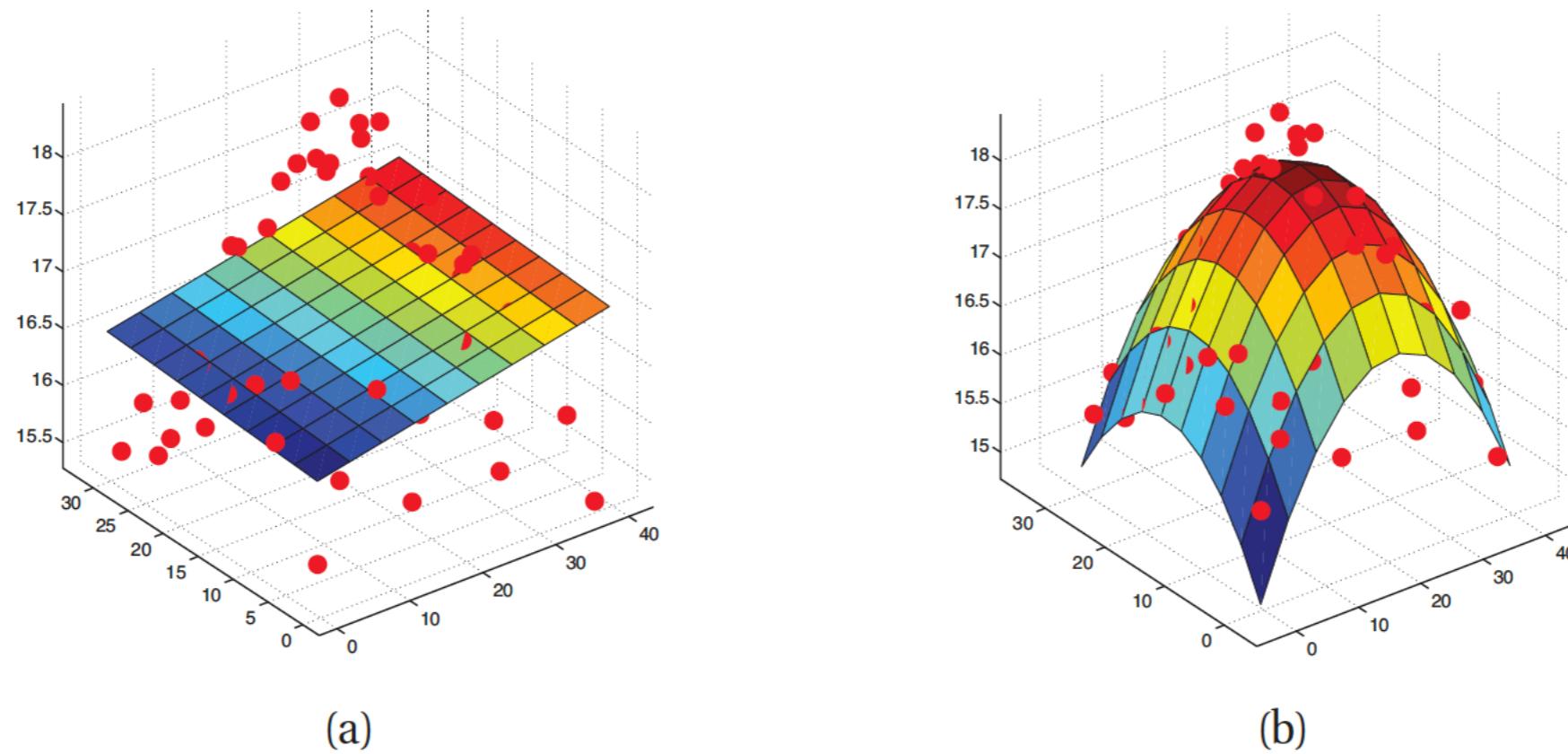
# Linear regression - Probability

Linear regression can be made to model non-linear relationships by replacing  $\mathbf{x}$  with some non-linear function of the inputs,  $\phi(\mathbf{x})$ . That is, we use

$$p(y|\mathbf{x}, \boldsymbol{\theta}) = \mathcal{N}(y|\mathbf{w}^T \phi(\mathbf{x}), \sigma^2) \quad (1.7)$$

This is known as **basis function expansion**. For example, Figure 1.18 illustrates the case where  $\phi(\mathbf{x}) = [1, x, x^2, \dots, x^d]$ , for  $d = 14$  and  $d = 20$ ; this is known as **polynomial regression**. We will consider other kinds of basis functions later in the book. In fact, many popular machine learning methods — such as support vector machines, neural networks, classification and regression trees, etc. — can be seen as just different ways of estimating basis functions from data, as we discuss in Chapters 14 and 16.

# Linear regression - Probability



**Figure 7.1** Linear regression applied to 2d data. Vertical axis is temperature, horizontal axes are location within a room. Data was collected by some remote sensing motes at Intel's lab in Berkeley, CA (data courtesy of Romain Thibaux). (a) The fitted plane has the form  $\hat{f}(\mathbf{x}) = w_0 + w_1x_1 + w_2x_2$ . (b) Temperature data is fitted with a quadratic of the form  $\hat{f}(\mathbf{x}) = w_0 + w_1x_1 + w_2x_2 + w_3x_1^2 + w_4x_2^2$ . Produced by `surfaceFitDemo`.

# Linear regression - MLE

## Maximum likelihood estimation (least squares)

A common way to estimate the parameters of a statistical model is to compute the MLE, which is defined as

$$\hat{\boldsymbol{\theta}} \triangleq \arg \max_{\boldsymbol{\theta}} \log p(\mathcal{D}|\boldsymbol{\theta}) \quad (7.4)$$

It is common to assume the training examples are independent and identically distributed, commonly abbreviated to **iid**. This means we can write the log-likelihood as follows:

$$\ell(\boldsymbol{\theta}) \triangleq \log p(\mathcal{D}|\boldsymbol{\theta}) = \sum_{i=1}^N \log p(y_i | \mathbf{x}_i, \boldsymbol{\theta}) \quad (7.5)$$

Instead of maximizing the log-likelihood, we can equivalently minimize the **negative log likelihood** or **NLL**:

$$\text{NLL}(\boldsymbol{\theta}) \triangleq - \sum_{i=1}^N \log p(y_i | \mathbf{x}_i, \boldsymbol{\theta}) \quad (7.6)$$

The NLL formulation is sometimes more convenient, since many optimization software packages are designed to find the minima of functions, rather than maxima.

# Linear regression - MLE

Now let us apply the method of MLE to the linear regression setting. Inserting the definition of the Gaussian into the above, we find that the log likelihood is given by

$$\ell(\boldsymbol{\theta}) = \sum_{i=1}^N \log \left[ \left( \frac{1}{2\pi\sigma^2} \right)^{\frac{1}{2}} \exp \left( -\frac{1}{2\sigma^2} (y_i - \mathbf{w}^T \mathbf{x}_i)^2 \right) \right] \quad (7.7)$$

$$= \frac{-1}{2\sigma^2} RSS(\mathbf{w}) - \frac{N}{2} \log(2\pi\sigma^2) \quad (7.8)$$

RSS stands for **residual sum of squares** and is defined by

$$RSS(\mathbf{w}) \triangleq \sum_{i=1}^N (y_i - \mathbf{w}^T \mathbf{x}_i)^2 \quad (7.9)$$

The RSS is also called the **sum of squared errors**, or SSE, and  $SSE/N$  is called the **mean squared error** or MSE. It can also be written as the square of the  $\ell_2$  **norm** of the vector of

Assumption: fixed sigma, independence responses

Minimize sum of squared error  $\Leftrightarrow$  Maximize log likelihood

# Gauss-Markov theorem

$$\hat{\beta} = (X^T X)^{-1} X^T y.$$

**Theorem 1 (Gaussian-Markov).** *Linear regression solution has the least variance among all unbiased linear estimators of  $\beta$ , i.e. if  $\beta_1 = \alpha_1^T y$  is an unbiased estimator of  $\beta$ , then  $\text{Var}(\beta_1) \geq \text{Var}(\beta)$  in the sense of matrix.*

*Proof.* Recall that the linear regression model is equivalent to

$$y = X\beta + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma^2).$$

We may write  $\alpha_1^T = (X^T X)^{-1} X^T + D$ .  $D$  is a  $(p+1) \times N$  matrix.

$$\begin{aligned} E(\beta_1) &= E(\alpha_1^T y) = E\left(\left((X^T X)^{-1} X^T + D\right)(X\beta + \epsilon)\right) \\ &= \left((X^T X)^{-1} X^T + D\right) X\beta = (1 + DX)\beta. \end{aligned}$$

Since  $\beta_1$  is an unbiased estimator of  $\beta$ , it follows that

$$DX = 0.$$

# Gauss-Markov theorem

Then

$$\begin{aligned}Var(\beta_1) &= Var(\alpha_1^T y) = \alpha_1^T Var(y) \alpha_1 \\&= \sigma^2 \left( (X^T X)^{-1} X^T + D \right) \left( (X^T X)^{-1} X^T + D \right)^T \\&= \sigma^2 \left( (X^T X)^{-1} X^T + D \right) \left( X (X X^T)^{-1} + D^T \right) \\&= \sigma^2 \left( (X X^T)^{-1} + (X^T X)^{-1} X^T D^T + D X (X X^T)^{-1} + D D^T \right) \\&= \sigma^2 (X X^T)^{-1} + \sigma^2 D D^T.\end{aligned}$$

$$\begin{aligned}Var(\hat{\beta}) &= Var(\alpha_1^T y) = (X^T X)^{-1} X^T Var(y) \left( (X^T X)^{-1} X^T \right)^T \\&= \sigma^2 (X^T X)^{-1} X^T X (X X^T)^{-1} = (X X^T)^{-1}.\end{aligned}$$

So

$Var(\beta_1) - Var(\hat{\beta})$  is a nonnegative definite matrix.

# Linear regression - Bayes

Assumption: fixed sigma, see 7.6.3 for general case

In linear regression, the likelihood is given by

$$p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \mu, \sigma^2) = \mathcal{N}(\mathbf{y}|\mu + \mathbf{X}\mathbf{w}, \sigma^2\mathbf{I}_N) \quad (7.52)$$

$$\propto \exp\left(-\frac{1}{2\sigma^2}(\mathbf{y} - \mu\mathbf{1}_N - \mathbf{X}\mathbf{w})^T(\mathbf{y} - \mu\mathbf{1}_N - \mathbf{X}\mathbf{w})\right) \quad (7.53)$$

form  $p(\mu) \propto 1$ , and then integrate it out to get **Exercise**

$$p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \sigma^2) \propto \exp\left(-\frac{1}{2\sigma^2}\|\mathbf{y} - \bar{\mathbf{y}}\mathbf{1}_N - \mathbf{X}\mathbf{w}\|_2^2\right) \quad (7.54)$$

where  $\bar{y} = \frac{1}{N} \sum_{i=1}^N y_i$  is the empirical mean of the output. For notational simplicity, we shall assume the output has been centered, and write  $\mathbf{y}$  for  $\mathbf{y} - \bar{y}\mathbf{1}_N$ .

# Linear regression - Bayes

## Exercise

$$p(y|\mathbf{x}, \boldsymbol{\theta}) = \mathcal{N}(y|w_0 + \mathbf{w}^\top \mathbf{x}, \sigma^2)$$

$$\hat{\mathbf{w}} = (\mathbf{X}_c^\top \mathbf{X}_c)^{-1} \mathbf{X}_c^\top \mathbf{y}_c = \left[ \sum_{i=1}^{N_{\mathcal{D}}} (\mathbf{x}_n - \bar{\mathbf{x}})(\mathbf{x}_n - \bar{\mathbf{x}})^\top \right]^{-1} \left[ \sum_{i=1}^{N_{\mathcal{D}}} (y_n - \bar{y})(\mathbf{x}_n - \bar{\mathbf{x}}) \right]$$

$$\hat{w}_0 = \frac{1}{N} \sum_n y_n - \frac{1}{N} \sum_n \mathbf{x}_n^\top \hat{\mathbf{w}} = \bar{y} - \bar{\mathbf{x}}^\top \hat{\mathbf{w}}$$

$$\mathbf{x}_n^c = \mathbf{x}_n - \bar{\mathbf{x}}$$

$$\mathbf{y}_c = \mathbf{y} - \bar{\mathbf{y}}$$

# Linear regression - Bayes

The conjugate prior to the above Gaussian likelihood is also a Gaussian, which we will denote by  $p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{w}_0, \mathbf{V}_0)$ . Using Bayes rule for Gaussians, Equation 4.125, the posterior is given by

$$p(\mathbf{w}|\mathbf{X}, \mathbf{y}, \sigma^2) \propto \mathcal{N}(\mathbf{w}|\mathbf{w}_0, \mathbf{V}_0)\mathcal{N}(\mathbf{y}|\mathbf{X}\mathbf{w}, \sigma^2\mathbf{I}_N) = \mathcal{N}(\mathbf{w}|\mathbf{w}_N, \mathbf{V}_N) \quad (7.55)$$

$$\mathbf{w}_N = \mathbf{V}_N \mathbf{V}_0^{-1} \mathbf{w}_0 + \frac{1}{\sigma^2} \mathbf{V}_N \mathbf{X}^T \mathbf{y} \quad (7.56)$$

$$\mathbf{V}_N^{-1} = \mathbf{V}_0^{-1} + \frac{1}{\sigma^2} \mathbf{X}^T \mathbf{X} \quad (7.57)$$

$$\mathbf{V}_N = \sigma^2 (\sigma^2 \mathbf{V}_0^{-1} + \mathbf{X}^T \mathbf{X})^{-1} \quad (7.58)$$

**Theorem 4.4.1** (Bayes rule for linear Gaussian systems). *Given a linear Gaussian system, as in Equation 4.124, the posterior  $p(\mathbf{x}|\mathbf{y})$  is given by the following:*

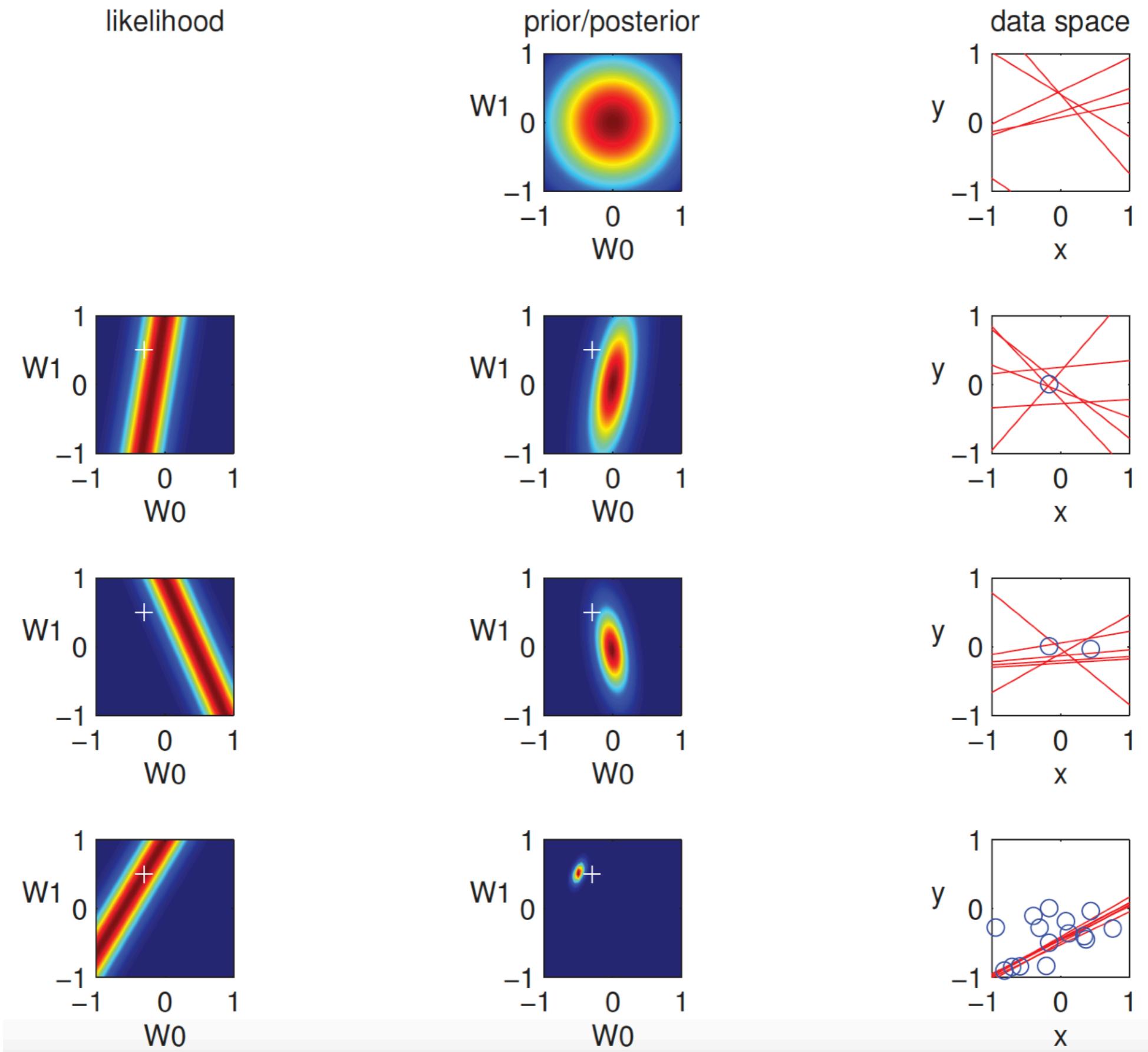
$$p(\mathbf{x}|\mathbf{y}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_{x|y}, \boldsymbol{\Sigma}_{x|y})$$

$$\boldsymbol{\Sigma}_{x|y}^{-1} = \boldsymbol{\Sigma}_x^{-1} + \mathbf{A}^T \boldsymbol{\Sigma}_y^{-1} \mathbf{A}$$

$$\boldsymbol{\mu}_{x|y} = \boldsymbol{\Sigma}_{x|y} [\mathbf{A}^T \boldsymbol{\Sigma}_y^{-1} (\mathbf{y} - \mathbf{b}) + \boldsymbol{\Sigma}_x^{-1} \boldsymbol{\mu}_x]$$

(4.125)

# Linear regression - Bayes



# Ridge regression

Ridge regression shrinks the regression coefficients by imposing a penalty on their size. The ridge coefficients minimize a penalized residual sum of squares,

$$\hat{\beta}^{\text{ridge}} = \underset{\beta}{\operatorname{argmin}} \left\{ \sum_{i=1}^N \left( y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2 + \lambda \sum_{j=1}^p \beta_j^2 \right\}. \quad (3.41)$$

Here  $\lambda \geq 0$  is a complexity parameter that controls the amount of shrink-

An equivalent way to write the ridge problem is

$$\hat{\beta}^{\text{ridge}} = \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^N \left( y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2,$$

$$\text{subject to } \sum_{j=1}^p \beta_j^2 \leq t,$$

# Ridge regression

replace  $x_{ij}$  with  $x_{ij} - \bar{x}_j$

Writing the criterion in (3.41) in matrix form,

$$\text{RSS}(\lambda) = (\mathbf{y} - \mathbf{X}\beta)^T(\mathbf{y} - \mathbf{X}\beta) + \lambda\beta^T\beta, \quad (3.43)$$

the ridge regression solutions are easily seen to be

$$\hat{\beta}^{\text{ridge}} = (\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I})^{-1}\mathbf{X}^T\mathbf{y}, \quad (3.44)$$

where  $\mathbf{I}$  is the  $p \times p$  identity matrix