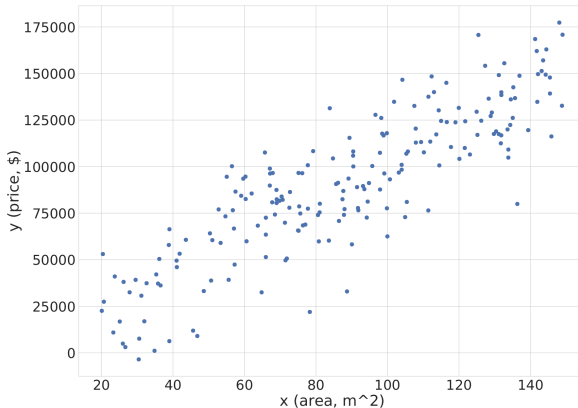


Section 1

Basic Linear Regression

Example: Housing price prediction

Given is a dataset $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$, of house areas x_i and corresponding prices y_i .



How do we estimate a price of a new house with area x_{new} ?

Regression problem

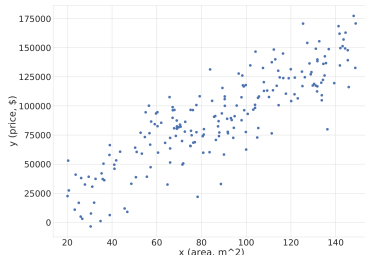
Given

- observations ¹
 $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}, \mathbf{x}_i \in \mathbb{R}^D$
- targets
 $\mathbf{y} = \{y_1, y_2, \dots, y_N\}, y_i \in \mathbb{R}$

Find

- Mapping $f(\cdot)$ from inputs to targets

$$y_i \approx f(\mathbf{x}_i)$$



¹A common way to represent the samples is as a **data matrix** $\mathbf{X} \in \mathbb{R}^{N \times D}$, where each row represents one sample.

Linear model

Target y is generated by a deterministic function f of \mathbf{x} plus noise

$$y_i = f(\mathbf{x}_i) + \epsilon_i, \quad \epsilon_i \sim \mathcal{N}(0, \beta^{-1}) \quad (1)$$

Let's choose $f(\mathbf{x})$ to be a linear function

$$f_{\mathbf{w}}(\mathbf{x}_i) = w_0 + w_1 x_{i1} + w_2 x_{i2} + \dots + w_D x_{iD} \quad (2)$$

$$= w_0 + \mathbf{w}^T \mathbf{x}_i \quad (3)$$

From now we will always assume that the bias term is absorbed into the \mathbf{x} vector

Absorbing the bias term

The linear function is given by

$$f_{\mathbf{w}}(\mathbf{x}) = w_0 + w_1x_1 + w_2x_2 + \dots + w_Dx_D \quad (4)$$

$$= w_0 + \mathbf{w}^T \mathbf{x} \quad (5)$$

Here w_0 is called **bias** or **offset** term. For simplicity, we can "absorb" it by prepending a 1 to the feature vector \mathbf{x} and respectively adding w_0 to the weight vector \mathbf{w} :

$$\tilde{\mathbf{x}} = (1, x_1, \dots, x_D)^T \quad \tilde{\mathbf{w}} = (w_0, w_1, \dots, w_D)^T$$

The function $f_{\mathbf{w}}$ can compactly be written as $f_{\mathbf{w}}(\mathbf{x}) = \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}$.

To unclutter the notation, we will assume the bias term is always absorbed and write \mathbf{w} and \mathbf{x} instead of $\tilde{\mathbf{w}}$ and $\tilde{\mathbf{x}}$.

Now, how do we choose the "best" \mathbf{w} that fits our data?

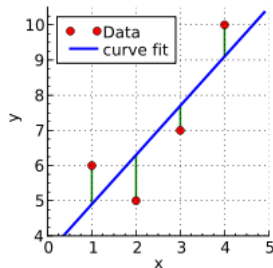
Error function

Error function gives a measure of "misfit" between our model (parametrized by \mathbf{w}) and observed data $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$.

Standard choice - **least squares** (LS) function

$$E_{\text{LS}}(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N (f_{\mathbf{w}}(\mathbf{x}_i) - y_i)^2 \quad (6)$$

$$= \frac{1}{2} \sum_{i=1}^N (\mathbf{w}^T \mathbf{x}_i - y_i)^2 \quad (7)$$



Factor $\frac{1}{2}$ is for later convenience

Objective

Find the optimal weight vector \mathbf{w}^* that minimizes the error

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} E_{\text{LS}}(\mathbf{w}) \quad (8)$$

$$= \arg \min_{\mathbf{w}} \frac{1}{2} \sum_{i=1}^N (\mathbf{x}_i^T \mathbf{w} - y_i)^2 \quad (9)$$

By stacking the observations \mathbf{x}_i as rows of the matrix $\mathbf{X} \in \mathbb{R}^{N \times D}$

$$= \arg \min_{\mathbf{w}} \frac{1}{2} (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y}) \quad (10)$$

Optimal solution

To find the minimum of the function $E(\mathbf{w})$, compute the gradient $\nabla_{\mathbf{w}} E(\mathbf{w})$:

$$\nabla_{\mathbf{w}} E_{\text{LS}}(\mathbf{w}) = \nabla_{\mathbf{w}} \frac{1}{2} (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y}) \quad (11)$$

$$= \nabla_{\mathbf{w}} \frac{1}{2} \left(\mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} - 2\mathbf{w}^T \mathbf{X}^T \mathbf{y} + \mathbf{y}^T \mathbf{y} \right) \quad (12)$$

$$= \mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{X}^T \mathbf{y} \quad (13)$$

See Equations (69), (81) from Matrix cookbook for details

Optimal solution

Now set the gradient to zero and solve for w to obtain the minimizer ²

$$X^T X w - X^T y \stackrel{!}{=} 0 \quad (14)$$

This leads to the so-called **normal equation** of the least squares problem

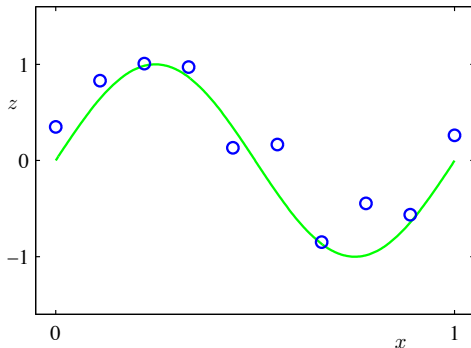
$$w^* = \underbrace{(X^T X)^{-1} X^T}_{=X^\dagger} y \quad (15)$$

X^\dagger is called **Moore-Penrose pseudo-inverse** of X (because for an invertible square matrix, $X^\dagger = X^{-1}$).

²Because Hessian $\nabla_w \nabla_w E(w)$ is positive (semi)definite \rightarrow see *Optimization*

Nonlinear dependency in data

What if the dependency between y and x is not linear?



Data generating process: $y_i = \sin(2\pi x_i) + \epsilon$, $\epsilon \sim \mathcal{N}(0, \beta^{-1})$

For this example assume that data dimensionality $D = 1$

Polynomials

Solution: Polynomials are universal function approximators, so we can define f as

$$f_{\mathbf{w}}(x) = w_0 + \sum_{j=1}^M w_j x^j \quad (16)$$

Or more generally

$$= w_0 + \sum_{j=1}^M w_j \phi_j(x) \quad (17)$$

Define $\phi_0 = 1$

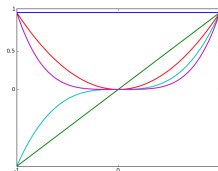
$$= \mathbf{w}^T \boldsymbol{\phi}(x) \quad (18)$$

The function f is still linear in \mathbf{w} (despite not being linear in x)!

Typical basis functions

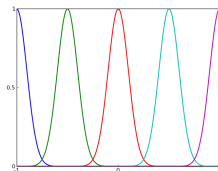
Polynomials

$$\phi_j(x) = x^j$$



Gaussian

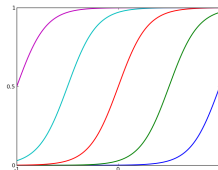
$$\phi_j(x) = e^{\frac{-(x-\mu_j)^2}{2s^2}}$$



Logistic Sigmoid

$$\phi_j(x) = \sigma\left(\frac{x-\mu_j}{s}\right),$$

where $\sigma(a) = \frac{1}{1+e^{-a}}$



Linear basis function model

Prediction for one sample

$$f_{\mathbf{w}}(\mathbf{x}) = w_0 + \sum_{j=1}^M w_j \phi_j(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}) \quad (19)$$

Using the same least squares error function as before

$$E_{\text{LS}}(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N (\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_i) - y_i)^2 \quad (20)$$

$$= \frac{1}{2} (\boldsymbol{\Phi} \mathbf{w} - \mathbf{y})^T (\boldsymbol{\Phi} \mathbf{w} - \mathbf{y}) \quad (21)$$

with

$$\boldsymbol{\Phi} = \begin{pmatrix} \phi_0(\mathbf{x}_1) & \phi_1(\mathbf{x}_1) & \dots & \phi_M(\mathbf{x}_1) \\ \phi_0(\mathbf{x}_2) & \phi_1(\mathbf{x}_2) & & \vdots \\ \vdots & \vdots & \ddots & \\ \phi_0(\mathbf{x}_N) & \phi_1(\mathbf{x}_N) & \dots & \phi_M(\mathbf{x}_N) \end{pmatrix} \in \mathbb{R}^{N \times (M+1)}$$

being the **design matrix** of $\boldsymbol{\phi}$.

Optimal solution

Recall Equation 10 - we have the same expression except that data matrix $\mathbf{X} \in \mathbb{R}^{N \times D}$ is replaced by design matrix $\Phi \in \mathbb{R}^{N \times (M+1)}$

$$E_{\text{LS}}(\mathbf{w}) = \frac{1}{2}(\Phi \mathbf{w} - \mathbf{y})^T (\Phi \mathbf{w} - \mathbf{y}) \quad (22)$$

This means that the optimal weights \mathbf{w}^* can be obtained in a similar way

$$\mathbf{w}^* = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{y} \quad (23)$$

$$= \Phi^\dagger \mathbf{y} \quad (24)$$

Compare this to Equation 15:

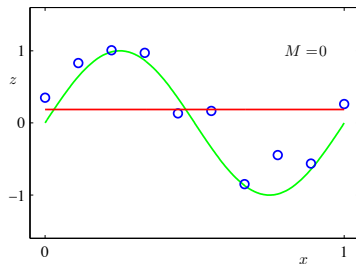
$$\mathbf{w}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad (25)$$

Choosing degree of the polynomial

How do we choose the degree of the polynomial M ?

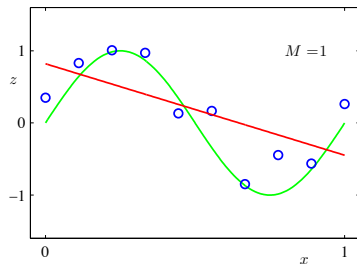
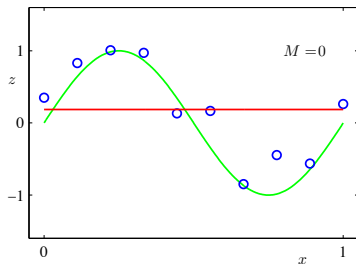
Choosing degree of the polynomial

How do we choose the degree of the polynomial M ?



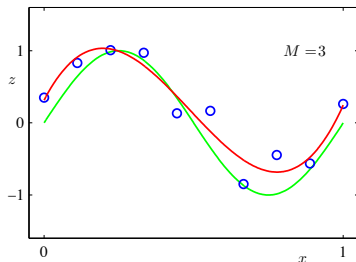
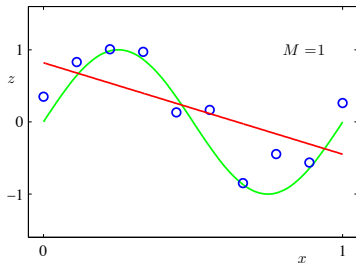
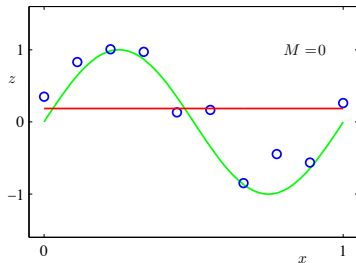
Choosing degree of the polynomial

How do we choose the degree of the polynomial M ?



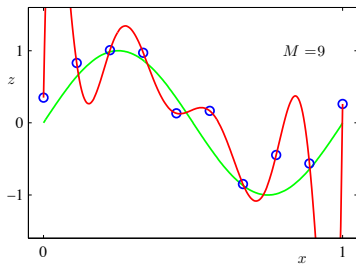
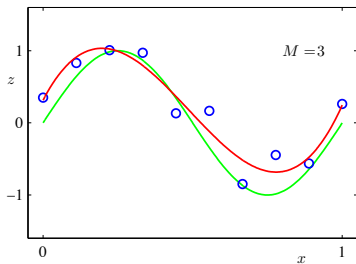
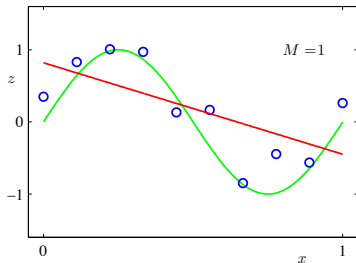
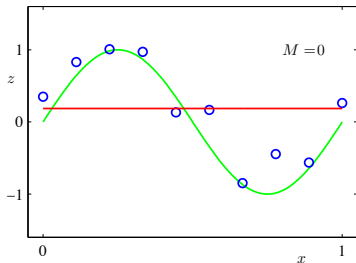
Choosing degree of the polynomial

How do we choose the degree of the polynomial M ?

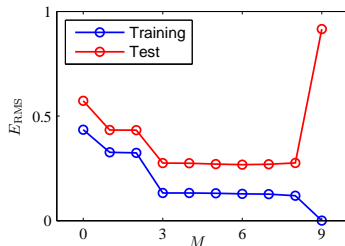
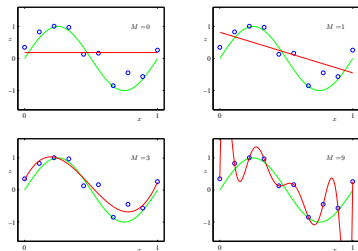


Choosing degree of the polynomial

How do we choose the degree of the polynomial M ?

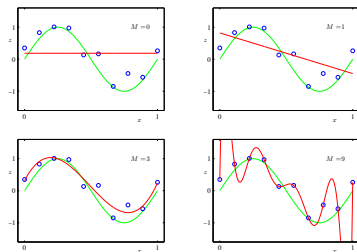


Choosing degree of the polynomial



One valid solution is to choose M using the standard train-validation split approach.

Choosing degree of the polynomial



	$M = 0$	$M = 1$	$M = 6$	$M = 9$
w_0^*	0.19	0.82	0.31	0.35
w_1^*		-1.27	7.99	232.37
w_2^*			-25.43	-5321.83
w_3^*			17.37	48568.31
w_4^*				-231639.30
w_5^*				640042.26
w_6^*				-1061800.52
w_7^*				1042400.18
w_8^*				-557682.99
w_9^*				125201.43

We also make another observation: overfitting occurs when the coefficients w become large.

What if we penalize large weights?

Controlling overfitting with regularization

Least squares loss function with L2 regularization
(also called ridge regression)

$$E_{\text{ridge}}(\mathbf{w}) = \frac{1}{2} \sum^N [\mathbf{w}^T \phi(\mathbf{x}_i) - y_i]^2 + \frac{\lambda}{2} \|\mathbf{w}\|_2^2 \quad (26)$$

Where,

- $\|\mathbf{w}\|_2^2 \equiv \mathbf{w}^T \mathbf{w} = w_0^2 + w_1^2 + w_2^2 + \dots + w_M^2$ - L2 norm
- λ - regularization strength

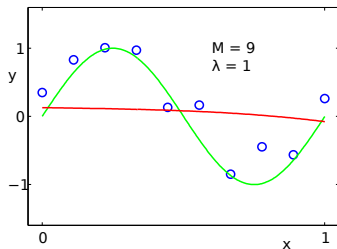
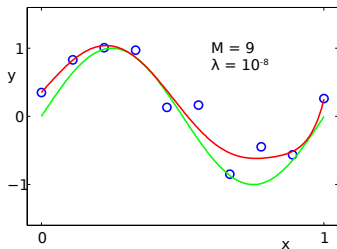
Controlling overfitting with regularization

Least squares loss function with L2 regularization
(also called **ridge regression**)

$$E_{\text{ridge}}(\mathbf{w}) = \frac{1}{2} \sum^N [\mathbf{w}^T \phi(\mathbf{x}_i) - y_i]^2 + \frac{\lambda}{2} \|\mathbf{w}\|_2^2 \quad (26)$$

Where,

- $\|\mathbf{w}\|_2^2 \equiv \mathbf{w}^T \mathbf{w} = w_0^2 + w_1^2 + w_2^2 + \dots + w_M^2$ - L2 norm
- λ - regularization strength



Larger regularization strength λ leads to smaller weights \mathbf{w}

Section 2

Probabilistic Linear Regression

Probabilistic formulation

Remember from our problem definition at the start of the lecture,

$$y_i = f_{\mathbf{w}}(\mathbf{x}_i) + \underbrace{\epsilon_i}_{\text{noise}}$$

Noise has zero-mean Gaussian distribution with a fixed precision $\beta = \frac{1}{\sigma^2}$

$$\epsilon_i \sim \mathcal{N}(0, \beta^{-1})$$

This implies that the distribution of the targets is

$$y_i \sim \mathcal{N}(f_{\mathbf{w}}(\mathbf{x}_i), \beta^{-1})$$

Remember: any function can be represented as $f_{\mathbf{w}}(\mathbf{x}_i) = \mathbf{w}^T \phi(\mathbf{x}_i)$

Maximum likelihood

Likelihood of a single sample

$$p(y_i | f_{\mathbf{w}}(\mathbf{x}_i), \beta) = \mathcal{N}(y_i | f_{\mathbf{w}}(\mathbf{x}_i), \beta^{-1}) \quad (27)$$

Assume that the samples are drawn i.i.d.

\implies likelihood of the entire dataset $\mathcal{D} = \{\mathbf{X}, \mathbf{y}\}$ is

$$p(\mathbf{y} | \mathbf{X}, \mathbf{w}, \beta) = \prod_{i=1}^N p(y_i | f_{\mathbf{w}}(\mathbf{x}_i), \beta) \quad (28)$$

We can now use the same approach we used in previous lecture - maximize the likelihood w.r.t. \mathbf{w} and β

$$\mathbf{w}_{\text{ML}}, \beta_{\text{ML}} = \arg \max_{\mathbf{w}, \beta} p(\mathbf{y} | \mathbf{X}, \mathbf{w}, \beta) \quad (29)$$

Maximum likelihood

Like in the coin flip example, we can make a few simplifications

$$\mathbf{w}_{\text{ML}}, \beta_{\text{ML}} = \arg \max_{\mathbf{w}, \beta} p(\mathbf{y} \mid \mathbf{X}, \mathbf{w}, \beta) \quad (30)$$

$$= \arg \max_{\mathbf{w}, \beta} \ln p(\mathbf{y} \mid \mathbf{X}, \mathbf{w}, \beta) \quad (31)$$

$$= \arg \min_{\mathbf{w}, \beta} -\ln p(\mathbf{y} \mid \mathbf{X}, \mathbf{w}, \beta) \quad (32)$$

Let's denote this quantity as **maximum likelihood error function** that we need to minimize

$$E_{\text{ML}}(\mathbf{w}, \beta) = -\ln p(\mathbf{y} \mid \mathbf{X}, \mathbf{w}, \beta) \quad (33)$$

Maximum likelihood

Simplify the error function

$$E_{\text{ML}}(\mathbf{w}, \beta) = -\ln \left[\prod_{i=1}^N \mathcal{N}(y_i \mid f_{\mathbf{w}}(\mathbf{x}_i), \beta^{-1}) \right] \quad (34)$$

$$= -\ln \left[\prod_{i=1}^N \sqrt{\frac{\beta}{2\pi}} \exp \left(-\frac{\beta}{2} (\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_i) - y_i)^2 \right) \right] \quad (35)$$

$$= -\sum_{i=1}^N \ln \left[\sqrt{\frac{\beta}{2\pi}} \exp \left(-\frac{\beta}{2} (\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_i) - y_i)^2 \right) \right] \quad (36)$$

$$= \frac{\beta}{2} \sum_{i=1}^N (\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_i) - y_i)^2 - \frac{N}{2} \ln \beta + \frac{N}{2} \ln 2\pi \quad (37)$$

Optimizing log-likelihood w.r.t. \mathbf{w}

$$\mathbf{w}_{\text{ML}} = \arg \min_{\mathbf{w}} E_{\text{ML}}(\mathbf{w}, \beta) \quad (38)$$

$$= \arg \min_{\mathbf{w}} \left[\frac{\beta}{2} \sum_{i=1}^N (\mathbf{w}^T \phi(\mathbf{x}_i) - y_i)^2 - \underbrace{\frac{N}{2} \ln \beta + \frac{N}{2} \ln 2\pi}_{= \text{const}} \right] \quad (39)$$

$$= \arg \min_{\mathbf{w}} \underbrace{\frac{1}{2} \sum_{i=1}^N (\mathbf{w}^T \phi(\mathbf{x}_i) - y_i)^2}_{\text{least squares error fn!}} \quad (40)$$

$$= \arg \min_{\mathbf{w}} E_{\text{LS}}(\mathbf{w}) \quad (41)$$

Optimizing log-likelihood w.r.t. \mathbf{w}

$$\mathbf{w}_{\text{ML}} = \arg \min_{\mathbf{w}} E_{\text{ML}}(\mathbf{w}, \beta) \quad (38)$$

$$= \arg \min_{\mathbf{w}} \left[\frac{\beta}{2} \sum_{i=1}^N (\mathbf{w}^T \phi(\mathbf{x}_i) - y_i)^2 - \underbrace{\frac{N}{2} \ln \beta + \frac{N}{2} \ln 2\pi}_{= \text{const}} \right] \quad (39)$$

$$= \arg \min_{\mathbf{w}} \underbrace{\frac{1}{2} \sum_{i=1}^N (\mathbf{w}^T \phi(\mathbf{x}_i) - y_i)^2}_{\text{least squares error fn!}} \quad (40)$$

$$= \arg \min_{\mathbf{w}} E_{\text{LS}}(\mathbf{w}) \quad (41)$$

Maximizing the likelihood is equivalent to minimizing the least squares error function!

$$\mathbf{w}_{\text{ML}} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{y} = \Phi^\dagger \mathbf{y} \quad (42)$$

Optimizing log-likelihood w.r.t. β

Plug in the estimate for \mathbf{w} and minimize w.r.t. β

$$\beta_{\text{ML}} = \arg \min_{\beta} E_{\text{ML}}(\mathbf{w}_{\text{ML}}, \beta) \quad (43)$$

$$= \arg \min_{\beta} \left[\frac{\beta}{2} \sum_{i=1}^N (\mathbf{w}_{\text{ML}}^T \phi(\mathbf{x}_i) - y_i)^2 - \frac{N}{2} \ln \beta + \frac{N}{2} \ln 2\pi \right] \quad (44)$$

Take derivative w.r.t. β and set it to zero

$$\frac{\partial}{\partial \beta} E_{\text{ML}}(\mathbf{w}_{\text{ML}}, \beta) = \frac{1}{2} \sum_{i=1}^N (\mathbf{w}_{\text{ML}}^T \phi(\mathbf{x}_i) - y_i)^2 - \frac{N}{2\beta} \stackrel{!}{=} 0 \quad (45)$$

Solving for β

$$\frac{1}{\beta_{\text{ML}}} = \frac{1}{N} \sum_{i=1}^N (\mathbf{w}_{\text{ML}}^T \phi(\mathbf{x}_i) - y_i)^2 \quad (46)$$

Predicting for new data

Recall, that

$$y \sim \mathcal{N}(f_{\mathbf{w}}(\mathbf{x}), \beta^{-1}) \quad (47)$$

Plugging in the \mathbf{w}_{ML} and β_{ML} into our likelihood we get a **predictive distribution** that allows us to make prediction \hat{y}_{new} for the new data \mathbf{x}_{new} .

$$p(\hat{y}_{\text{new}} \mid \mathbf{x}_{\text{new}}, \mathbf{w}_{\text{ML}}, \beta_{\text{ML}}) = \mathcal{N}(\hat{y}_{\text{new}} \mid \mathbf{w}_{\text{ML}}^T \boldsymbol{\phi}(\mathbf{x}_{\text{new}}), \beta_{\text{ML}}^{-1}) \quad (48)$$

Posterior distribution

Recall from the Lecture 3, that ML leads to overfitting (especially, when little training data is available).

Solution - consider the **posterior distribution** instead

$$p(\mathbf{w} \mid \mathbf{X}, \mathbf{y}, \beta, \cdot) = \frac{\overbrace{p(\mathbf{y} \mid \mathbf{X}, \mathbf{w}, \beta)}^{\text{likelihood}} \cdot \overbrace{p(\mathbf{w} \mid \cdot)}^{\text{prior}}}{\underbrace{p(\mathbf{X}, \mathbf{y})}_{\text{normalizing constant}}} \quad (49)$$

$$\propto p(\mathbf{y} \mid \mathbf{X}, \mathbf{w}, \beta) \cdot p(\mathbf{w} \mid \cdot) \quad (50)$$

Precision $\beta = 1/\sigma^2$ is treated as a known parameter to simplify the calculations.

Posterior distribution

Recall from the Lecture 3, that ML leads to overfitting (especially, when little training data is available).

Solution - consider the **posterior distribution** instead

$$p(\mathbf{w} \mid \mathbf{X}, \mathbf{y}, \beta, \cdot) = \frac{\overbrace{p(\mathbf{y} \mid \mathbf{X}, \mathbf{w}, \beta)}^{\text{likelihood}} \cdot \overbrace{p(\mathbf{w} \mid \cdot)}^{\text{prior}}}{\underbrace{p(\mathbf{X}, \mathbf{y})}_{\text{normalizing constant}}} \quad (49)$$

$$\propto p(\mathbf{y} \mid \mathbf{X}, \mathbf{w}, \beta) \cdot p(\mathbf{w} \mid \cdot) \quad (50)$$

Connection to the coin flip example

	train data	likelihood	prior	posterior
coin:	$\mathcal{D} = \mathbf{X}$	$p(\mathcal{D} \mid \theta)$	$p(\theta \mid a, b)$	$p(\theta \mid \mathcal{D})$
regr.:	$\mathcal{D} = \{\mathbf{X}, \mathbf{y}\}$	$p(\mathbf{y} \mid \mathbf{X}, \mathbf{w}, \beta)$	$p(\mathbf{w} \mid \cdot)$	$p(\mathbf{w} \mid \mathbf{X}, \mathbf{y}, \beta, \cdot)$

How do we choose the prior $p(\mathbf{w} \mid \cdot)$?

Precision $\beta = 1/\sigma^2$ is treated as a known parameter to simplify the calculations.

Prior for \mathbf{w}

We set the prior over \mathbf{w} to an isotropic multivariate normal distribution with zero mean

$$p(\mathbf{w} \mid \alpha) = \mathcal{N}(\mathbf{w} \mid \mathbf{0}, \alpha^{-1} \mathbf{I}) = \left(\frac{\alpha}{2\pi} \right)^{\frac{M}{2}} \exp \left(-\frac{\alpha}{2} \mathbf{w}^T \mathbf{w} \right) \quad (51)$$

where,

α - precision of the distribution

M - number of elements in the vector \mathbf{w}

Motivation:

- Higher probability is assigned to small values of \mathbf{w}
 \implies prevents overfitting (recall slide 21)
- Likelihood is also Gaussian - simplified calculations

Maximum a posteriori (MAP)

We are looking for \mathbf{w} that corresponds to the mode of the posterior

$$\mathbf{w}_{\text{MAP}} = \arg \max_{\mathbf{w}} p(\mathbf{w} \mid \mathbf{X}, \mathbf{y}, \alpha, \beta) \quad (52)$$

$$= \arg \max_{\mathbf{w}} \ln p(\mathbf{y} \mid \mathbf{X}, \mathbf{w}, \beta) + \ln p(\mathbf{w} \mid \alpha) - \underbrace{\ln p(\mathbf{X}, \mathbf{y})}_{=\text{const}} \quad (53)$$

$$= \arg \min_{\mathbf{w}} -\ln p(\mathbf{y} \mid \mathbf{X}, \mathbf{w}, \beta) - \ln p(\mathbf{w} \mid \alpha) \quad (54)$$

Similar to ML, define the MAP error function as negative log-posterior

$$E_{\text{MAP}}(\mathbf{w}) = -\ln p(\mathbf{w} \mid \mathbf{X}, \mathbf{y}, \alpha, \beta) \quad (55)$$

$$= -\ln p(\mathbf{y} \mid \mathbf{X}, \mathbf{w}, \beta) - \ln p(\mathbf{w} \mid \alpha) + \text{const} \quad (56)$$

We ignore the constant terms in the error function, as they are independent of \mathbf{w}

MAP error function

Simplify the error function

$$\begin{aligned}E_{MAP} &= -\ln p(\mathbf{y} \mid \mathbf{X}, \mathbf{w}, \beta) - \ln p(\mathbf{w} \mid \alpha) \\&= \frac{\beta}{2} \sum_{i=1}^N (\mathbf{w}^T \phi(\mathbf{x}_i) - y_i)^2 - \frac{N}{2} \ln \beta + \frac{N}{2} \ln 2\pi \\&\quad - \ln \left(\frac{\alpha}{2\pi} \right)^{\frac{M}{2}} + \frac{\alpha}{2} \mathbf{w}^T \mathbf{w} \\&= \frac{\beta}{2} \sum_{i=1}^N (\mathbf{w}^T \phi(\mathbf{x}_i) - y_i)^2 + \frac{\alpha}{2} \|\mathbf{w}\|_2^2 + \text{const} \\&= \underbrace{\frac{1}{2} \sum_{i=1}^N (\mathbf{w}^T \phi(\mathbf{x}_i) - y_i)^2 + \frac{\lambda}{2} \|\mathbf{w}\|_2^2}_{\text{ridge regression error fn!}} + \text{const} \quad \text{where } \lambda = \frac{\alpha}{\beta} \\&= E_{\text{ridge}}(\mathbf{w}) + \text{const}\end{aligned} \tag{57}$$

MAP estimation with Gaussian prior is equivalent to ridge regression!

Predicting for new data

Recall, that

$$y \sim \mathcal{N}(f_{\mathbf{w}}(\mathbf{x}), \beta^{-1}) \quad (58)$$

Plugging in the \mathbf{w}_{MAP} into our likelihood we get a **predictive distribution** that lets us make prediction \hat{y}_{new} for the new data \mathbf{x}_{new} .

$$p(\hat{y}_{\text{new}} \mid \mathbf{x}_{\text{new}}, \mathbf{w}_{\text{MAP}}, \beta) = \mathcal{N}(\hat{y}_{\text{new}} \mid \mathbf{w}_{\text{MAP}}^T \boldsymbol{\phi}(\mathbf{x}_{\text{new}}), \beta^{-1}) \quad (59)$$

Recall, that we assume β to be known a priori (for simplified calculations).