

Machine Learning

Lecture 6: Linear Regression

Prof. Dr. Aleksandar Bojchevski

02.05.23

Outline

Simple linear regression

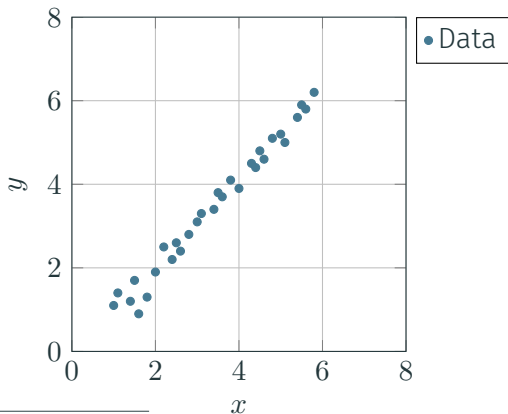
Non-linear relationships

Probabilistic Interpretation

Regression problem

Given a dataset $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$ of observations, predict y_{new} for a x_{new} ?

Given features $\mathbf{x}_i \in \mathbb{R}^D$ and labels $y_i \in \mathbb{R}$, find a mapping f such that $y_i \approx f(\mathbf{x}_i)$

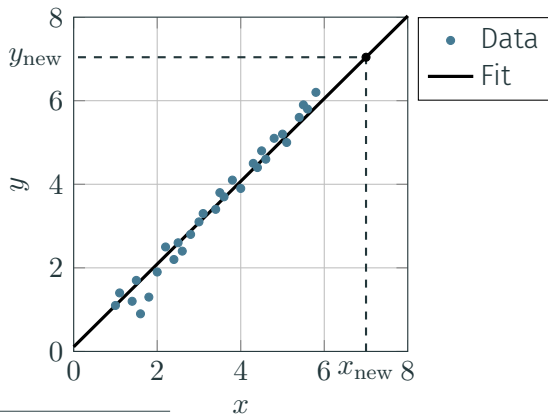


On the figure $x_i \in \mathbb{R}$ is one dimensional.

Regression problem

Given a dataset $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$ of observations, predict y_{new} for a x_{new} ?

Given features $\mathbf{x}_i \in \mathbb{R}^D$ and labels $y_i \in \mathbb{R}$, find a mapping f such that $y_i \approx f(\mathbf{x}_i)$



On the figure $x_i \in \mathbb{R}$ is one dimensional.

Linear model

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

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

We choose $f(\mathbf{x})$ to be a linear function

$$\begin{aligned} f_{\mathbf{w}}(\mathbf{x}_i) &= w_0 + w_1 x_{i1} + w_2 x_{i2} + \dots + w_D x_{iD} \\ &= w_0 + \mathbf{w}^T \mathbf{x}_i \end{aligned}$$

Here w_0 is called the bias or offset term.

Absorbing the bias term

The linear function is given by

$$\begin{aligned}f_{\mathbf{w}}(\mathbf{x}) &= w_0 + w_1x_1 + w_2x_2 + \dots + w_Dx_D \\ &= w_0 + \mathbf{w}^T \mathbf{x}\end{aligned}$$

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 $f_{\mathbf{w}}(\mathbf{x}) = \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}$.

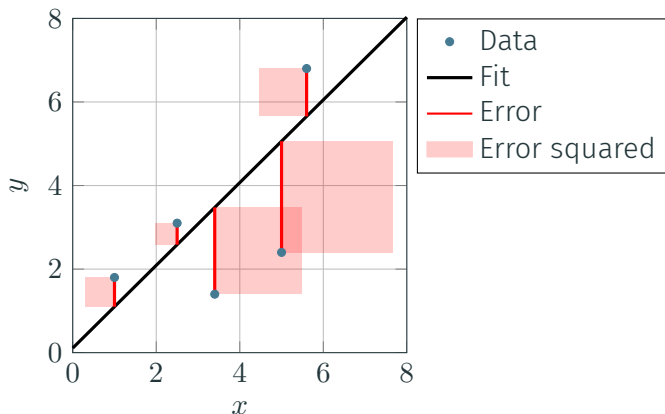
To unclutter the notation, we will assume the bias term is always absorbed and write \mathbf{w} and \mathbf{x} instead of $\tilde{\mathbf{w}}$ and $\tilde{\mathbf{x}}$.

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

The least squares loss

Standard choice to measure the error between our model (parametrized by \mathbf{w}) and the observed data $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$ is the **least squares** loss.

$$\begin{aligned}\mathcal{L}_{\text{LS}}(\mathbf{w}) &= \frac{1}{2} \sum_{i=1}^N (\hat{y}_i - y_i)^2 \\ &= \frac{1}{2} \sum_{i=1}^N (f_{\mathbf{w}}(\mathbf{x}_i) - y_i)^2 \\ &= \frac{1}{2} \sum_{i=1}^N (\mathbf{w}^T \mathbf{x}_i - y_i)^2\end{aligned}$$



The factor $\frac{1}{2}$ is for later convenience.

Minimizing the loss

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

$$\begin{aligned}\mathbf{w}^* &= \arg \min_{\mathbf{w}} \mathcal{L}(\mathbf{w}) \\ &= \arg \min_{\mathbf{w}} \frac{1}{2} \sum_{i=1}^N (\mathbf{x}_i^T \mathbf{w} - y_i)^2\end{aligned}$$

By stacking the observations \mathbf{x}_i as rows of a matrix $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)^T \in \mathbb{R}^{N \times D}$ and the targets y_i into a vector $\mathbf{y} = (y_1, \dots, y_N) \in \mathbb{R}^N$.

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

We omit the LS subscript in $\mathcal{L}_{\text{LS}}(\mathbf{w})$ when it is clear from the context.

To find the minimum of the loss $\mathcal{L}(\mathbf{w})$, compute the gradient $\nabla_{\mathbf{w}}\mathcal{L}(\mathbf{w})$:

$$\begin{aligned}\nabla_{\mathbf{w}}\mathcal{L}(\mathbf{w}) &= \nabla_{\mathbf{w}}\frac{1}{2}(\mathbf{X}\mathbf{w} - \mathbf{y})^T(\mathbf{X}\mathbf{w} - \mathbf{y}) \\ &= \nabla_{\mathbf{w}}\frac{1}{2}(\mathbf{w}^T\mathbf{X}^T\mathbf{X}\mathbf{w} - 2\mathbf{w}^T\mathbf{X}^T\mathbf{y} + \mathbf{y}^T\mathbf{y}) \\ &= \mathbf{X}^T\mathbf{X}\mathbf{w} - \mathbf{X}^T\mathbf{y}\end{aligned}$$

See Equations (69), (81) from The Matrix Cookbook for details.

Optimal solution

Now set the gradient to zero and solve for \mathbf{w} to obtain the minimizer¹

$$\begin{aligned}\mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{X}^T \mathbf{y} &\stackrel{!}{=} 0 \\ \mathbf{X}^T \mathbf{X} \mathbf{w} &= \mathbf{X}^T \mathbf{y}\end{aligned}$$

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

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

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

¹Because the Hessian $\nabla_{\mathbf{w}} \nabla_{\mathbf{w}} E(\mathbf{w}) = \mathbf{X}^T \mathbf{X}$ is positive (semi)definite.

Geometric interpretation

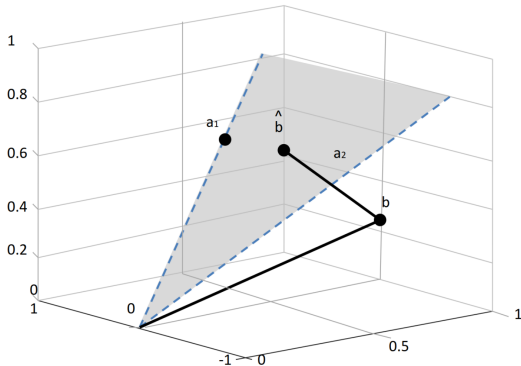
Assume more observations than unknowns ($N > D$) – an overdetermined system.

Example $\mathbf{Ax} = \mathbf{b}$ for $\mathbf{A} \in \mathbb{R}^{3 \times 2}$.

3 equations and 2 unknowns.

The columns of \mathbf{A} , \mathbf{a}_1 and \mathbf{a}_2 , define a 2D linear subspace embedded in \mathbb{R}^3 .

$\hat{\mathbf{b}}$ is the orthogonal projection of the target \mathbf{b} onto the linear subspace.



In the special case that $\mathbf{X} = \mathbf{x}$ is a column vector, the orthogonal projection of \mathbf{y} onto the line \mathbf{x} becomes $\mathbf{x} \frac{\mathbf{x}^T \mathbf{y}}{\mathbf{x}^T \mathbf{x}}$.

Algorithmic issues

The closed form ordinary least squares (OLS) solution is $\mathbf{w}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$.

Possible numerical issues if $\mathbf{X}^T \mathbf{X}$ is ill-conditioned or singular leading.

Alternatives:

- Compute the pseudo-inverse using the SVD (default in sklearn)

- If \mathbf{X} is tall and skinny ($N \gg D$) we can use QR decomposition²

- Iterative solvers such as conjugate gradient methods and GMRES³

²Here $\mathbf{w}^* = \mathbf{R}^{-1} \mathbf{Q}^T \mathbf{y}$, where $\mathbf{X} = \mathbf{Q}\mathbf{R}$ and $\mathbf{Q}^T \mathbf{Q} = \mathbf{I}$ and \mathbf{R} is upper triangular.

³Suitable when \mathbf{X} is sparse or structured since they only need the ability to perform matrix-vector multiplications.

Outline

Simple linear regression

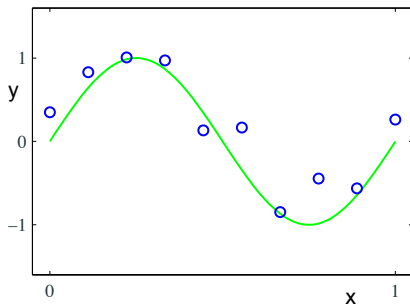
Non-linear relationships

Probabilistic Interpretation

Nonlinear dependency in data

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

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



For this example assume that the data dimensionality is $D = 1$.

Polynomials

Solution: Polynomials are universal function approximators.

So for a 1-dimensional x we can define f as

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

or more generally for any basis function $\phi_j(x)$

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

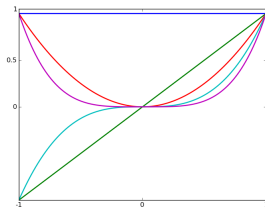
Absorbing the bias term by defining $\phi_0 = 1$

$$f_{\mathbf{w}}(x) = \mathbf{w}^T \boldsymbol{\phi}(x) = \mathbf{w}^T [\phi_0(x), \phi_1(x), \dots, \phi_M(x)]$$

The function f is still linear in \mathbf{w} despite being potentially non-linear in x !

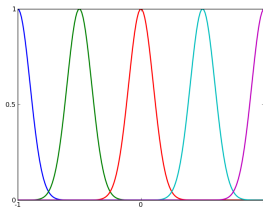
Typical basis functions

Polynomials



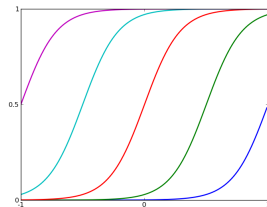
$$\phi_j(x) = x^j$$

Gaussian



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

Logistic Sigmoid



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

where $\sigma(a) = \frac{1}{1+\exp^{-a}}$ is the sigmoid, and μ_j and s are hyperparameters.

Linear basis function model

For D -dimensional data \mathbf{x} , $\phi_j : \mathbb{R}^D \rightarrow \mathbb{R}$ and we have

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

Using the same least squares error function as before

$$\mathcal{L}(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N (\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_i) - y_i)^2 = \frac{1}{2} (\boldsymbol{\Phi} \mathbf{w} - \mathbf{y})^T (\boldsymbol{\Phi} \mathbf{w} - \mathbf{y})$$

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

$\boldsymbol{\Phi}$ is called the **design matrix** of $\boldsymbol{\phi}$.

Optimal solution

Recall least squares loss for the original feature matrix \mathbf{X}

$$\mathcal{L}(\mathbf{w}) = \frac{1}{2}(\mathbf{X}\mathbf{w} - \mathbf{y})^T(\mathbf{X}\mathbf{w} - \mathbf{y})$$

and compare it to the expression we found with the design matrix $\Phi \in \mathbb{R}^{N \times (M+1)}$

$$\mathcal{L}(\mathbf{w}) = \frac{1}{2}(\Phi\mathbf{w} - \mathbf{y})^T(\Phi\mathbf{w} - \mathbf{y}).$$

The optimal weights \mathbf{w}^* can be obtained in the same way

$$\mathbf{w}^* = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{y} = \Phi^\dagger \mathbf{y} \tag{2}$$

Compare this to Equation 1: $\mathbf{w}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = \mathbf{X}^\dagger \mathbf{y}$

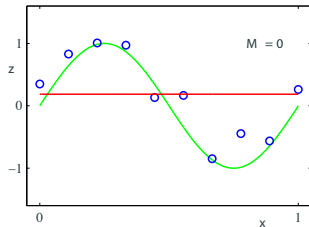
Only thing we've done is transformed our original features \mathbf{X} into new features Φ .

Choosing the degree of the polynomial

How do we choose the degree of the polynomial M ?

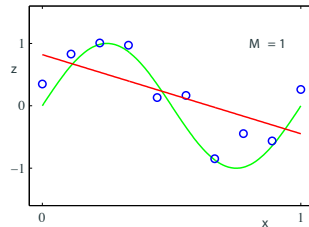
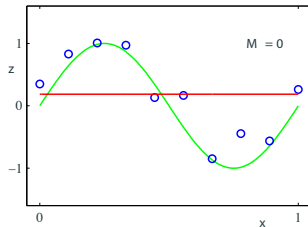
Choosing the degree of the polynomial

How do we choose the degree of the polynomial M ?



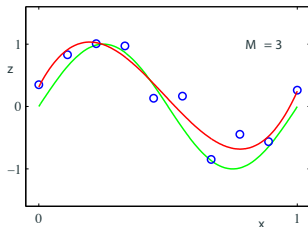
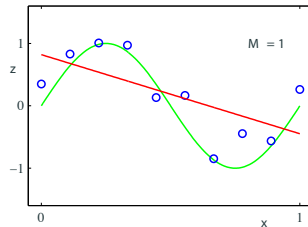
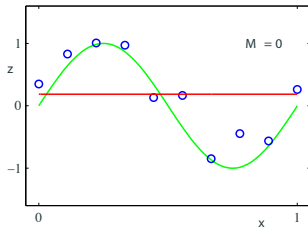
Choosing the degree of the polynomial

How do we choose the degree of the polynomial M ?



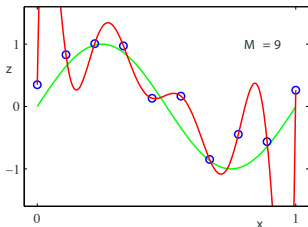
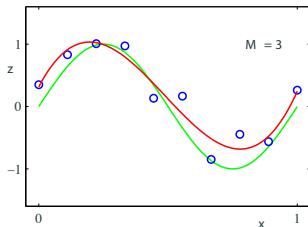
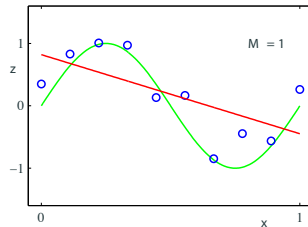
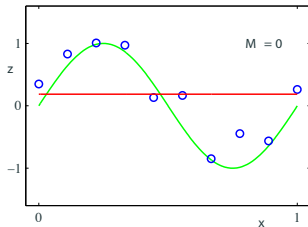
Choosing the degree of the polynomial

How do we choose the degree of the polynomial M ?



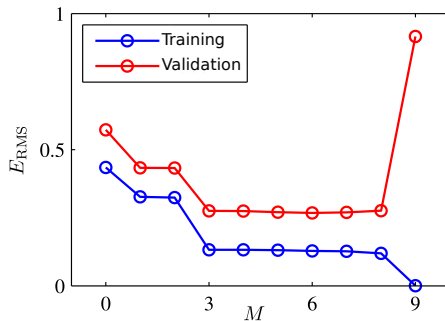
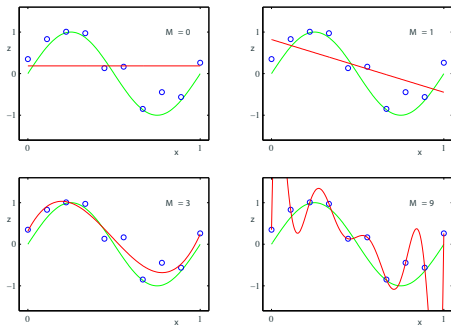
Choosing the degree of the polynomial

How do we choose the degree of the polynomial M ?



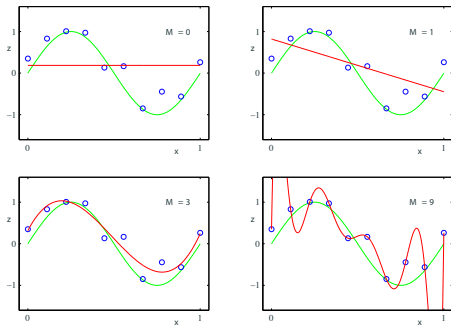
Choosing the degree of the polynomial

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



Manifestation of overfitting

Observation: overfitting occurs when the coefficients w become large.



| | $M = 0$ | $M = 1$ | $M = 3$ | $M = 9$ |
|---------|---------|---------|---------|-------------|
| 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 |

What if we penalize large weights?

Controlling overfitting with regularization

Least squares loss with **L2 regularization**, also called **ridge regression**

$$\mathcal{L}_{\text{ridge}}(\mathbf{w}) = \frac{1}{2} \sum^N [\mathbf{w}^T \phi(\mathbf{x}_i) - y_i]^2 + \frac{\lambda}{2} \|\mathbf{w}\|_2^2 \quad (3)$$

where $\|\mathbf{w}\|_2^2 \equiv \mathbf{w}^T \mathbf{w} = w_0^2 + w_1^2 + w_2^2 + \cdots + w_M^2$ is the squared L2 norm of \mathbf{w} and λ is the regularization strength.

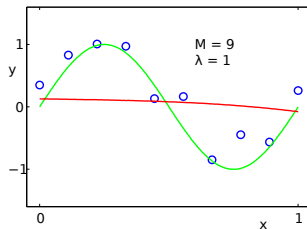
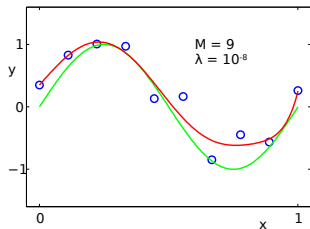
Controlling overfitting with regularization

Least squares loss with **L2 regularization**, also called **ridge regression**

$$\mathcal{L}_{\text{ridge}}(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N [\mathbf{w}^T \phi(\mathbf{x}_i) - y_i]^2 + \frac{\lambda}{2} \|\mathbf{w}\|_2^2 \quad (3)$$

where $\|\mathbf{w}\|_2^2 \equiv \mathbf{w}^T \mathbf{w} = w_0^2 + w_1^2 + w_2^2 + \dots + w_M^2$ is the squared L2 norm of \mathbf{w} and λ is the regularization strength.

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



Increasing the degree of the polynomial M

Increasing the regularization strength λ

Increasing the degree of the polynomial M

decreases bias – the function is more flexible

Increasing the regularization strength λ

Bias-variance trade-off

Increasing the degree of the polynomial M

- decreases bias – the function is more flexible

- increases variance – we can more easily fit the noise

Increasing the regularization strength λ

Bias-variance trade-off

Increasing the degree of the polynomial M

- decreases bias – the function is more flexible

- increases variance – we can more easily fit the noise

Increasing the regularization strength λ

- increases bias

Bias-variance trade-off

Increasing the degree of the polynomial M

- decreases bias – the function is more flexible

- increases variance – we can more easily fit the noise

Increasing the regularization strength λ

- increases bias

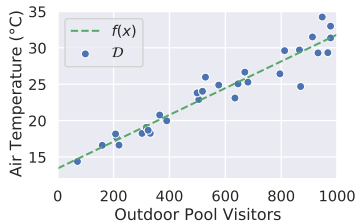
- decreases variance

Interpretation and correlation vs. causation

The weight w_d can be interpreted as the strength of the (linear) relationship between feature x_d and y .

Least squares fit $f(x) = 0.018x + 13.43$.

A weight of 0.018 shows a strong correlation (considering the different scales).



Correlation \nRightarrow causation: Increasing pool visitors doesn't increase temperature.

If you normalize the data to handle the different scales of x and y (as we often do in practice) we find a weight of about 1.

Simple linear regression

Non-linear relationships

Probabilistic Interpretation

Maximum likelihood

Assuming Gaussian likelihood (with a fixed precision β) for a single sample

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

and assuming that the samples are drawn independently, the likelihood of the entire dataset $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N = \{\mathbf{X}, \mathbf{y}\}$ is

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

We use the same approach as previously – maximize the likelihood w.r.t. \mathbf{w} and β

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

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

$$\begin{aligned}\mathbf{w}_{\text{ML}}, \beta_{\text{ML}} &= \arg \max_{\mathbf{w}, \beta} p(\mathbf{y} \mid \mathbf{X}, \mathbf{w}, \beta) \\ &= \arg \max_{\mathbf{w}, \beta} \ln p(\mathbf{y} \mid \mathbf{X}, \mathbf{w}, \beta) \\ &= \arg \min_{\mathbf{w}, \beta} -\ln p(\mathbf{y} \mid \mathbf{X}, \mathbf{w}, \beta)\end{aligned}$$

We need to minimize the negative log-likelihood which is our loss

$$\mathcal{L}_{\text{NLL}}(\mathbf{w}, \beta) = -\ln p(\mathbf{y} \mid \mathbf{X}, \mathbf{w}, \beta)$$

$$\begin{aligned}\mathcal{L}_{\text{NLL}}(\mathbf{w}, \beta) &= -\ln \left[\prod_{i=1}^N \mathcal{N}(y_i \mid f_{\mathbf{w}}(\mathbf{x}_i), \beta^{-1}) \right] \\&= -\ln \left[\prod_{i=1}^N \sqrt{\frac{\beta}{2\pi}} \exp \left(-\frac{\beta}{2} (\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_i) - y_i)^2 \right) \right] \\&= -\sum_{i=1}^N \ln \left[\sqrt{\frac{\beta}{2\pi}} \exp \left(-\frac{\beta}{2} (\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_i) - y_i)^2 \right) \right] \\&= \frac{\beta}{2} \sum_{i=1}^N (\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_i) - y_i)^2 - \frac{N}{2} \ln \beta + \frac{N}{2} \ln 2\pi\end{aligned}$$

Minimizing the negative log-likelihood w.r.t. \mathbf{w}

$$\begin{aligned}\mathbf{w}_{\text{ML}} &= \arg \min_{\mathbf{w}} \mathcal{L}_{\text{NLL}}(\mathbf{w}, \beta) \\&= \arg \min_{\mathbf{w}} \left[\frac{\beta}{2} \sum_{i=1}^N (\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_i) - y_i)^2 - \underbrace{\frac{N}{2} \ln \beta + \frac{N}{2} \ln 2\pi}_{=\text{const}} \right] \\&= \arg \min_{\mathbf{w}} \underbrace{\frac{1}{2} \sum_{i=1}^N (\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_i) - y_i)^2}_{\text{least squares loss}} \\&= \arg \min_{\mathbf{w}} \mathcal{L}_{\text{LS}}(\mathbf{w})\end{aligned}$$

Maximizing the likelihood is equivalent to minimizing the least squares loss!

Thus, we have $\mathbf{w}^* = \mathbf{w}_{\text{LS}} = \mathbf{w}_{\text{ML}} = (\boldsymbol{\Phi}^T \boldsymbol{\Phi})^{-1} \boldsymbol{\Phi}^T \mathbf{y} = \boldsymbol{\Phi}^\dagger \mathbf{y}$

Minimizing the negative log-likelihood w.r.t. β

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

$$\begin{aligned}\beta_{\text{ML}} &= \arg \min_{\beta} E_{\text{ML}}(\mathbf{w}_{\text{ML}}, \beta) \\ &= \arg \min_{\beta} \left[\frac{\beta}{2} \sum_{i=1}^N (\mathbf{w}_{\text{ML}}^T \phi(\mathbf{x}_i) - y_i)^2 - \frac{N}{2} \ln \beta + \frac{N}{2} \ln 2\pi \right]\end{aligned}$$

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

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

Solving for β

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

Posterior distribution

Since MLE can overfit with little data consider the **posterior distribution** instead.

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

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

Posterior distribution

Since MLE can overfit with little data consider the **posterior distribution** instead.

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

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.: | $\mathcal{D} = \{\mathbf{X}, \mathbf{y}\}$ | $p(\mathbf{y} \mid \mathbf{X}, \mathbf{w}, \beta)$ | $p(\mathbf{w} \mid \cdot)$ | $p(\mathbf{w} \mid \mathbf{X}, \mathbf{y}, \beta, \cdot)$ |

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

Choosing the prior for \mathbf{w}

Isotropic multivariate normal distribution with zero mean as prior for \mathbf{w}

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

where α is the precision, and M is the dimension.

Motivation:

Higher probability is assigned to small values of $\mathbf{w} \implies$ prevents overfitting
(recall slide 20)

Likelihood is also Gaussian - simplified calculations

Maximum a posteriori (MAP)

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

$$\begin{aligned}\mathbf{w}_{\text{MAP}} &= \arg \max_{\mathbf{w}} p(\mathbf{w} \mid \mathbf{X}, \mathbf{y}, \alpha, \beta) \\ &= \arg \max_{\mathbf{w}} \ln p(\mathbf{y} \mid \mathbf{X}, \mathbf{w}, \beta) + \ln p(\mathbf{w} \mid \alpha) - \underbrace{\ln p(\mathbf{y} \mid \mathbf{X}, \beta, \alpha)}_{=\text{const}} \\ &= \arg \min_{\mathbf{w}} -\ln p(\mathbf{y} \mid \mathbf{X}, \mathbf{w}, \beta) - \ln p(\mathbf{w} \mid \alpha)\end{aligned}$$

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

$$\begin{aligned}\mathcal{L}_{\text{MAP}}(\mathbf{w}) &= -\ln p(\mathbf{w} \mid \mathbf{X}, \mathbf{y}, \alpha, \beta) \\ &= -\ln p(\mathbf{y} \mid \mathbf{X}, \mathbf{w}, \beta) - \ln p(\mathbf{w} \mid \alpha) + \text{const}\end{aligned}$$

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

Simplify the MAP loss function

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

MAP estimation with Gaussian prior is equivalent to ridge regression!

Full Bayesian approach

Full posterior $p(\mathbf{w} \mid \mathcal{D}) \propto p(\mathbf{y} \mid \mathbf{X}, \mathbf{w}, \beta) p(\mathbf{w} \mid \alpha)$ instead of a point estimate \mathbf{w}_{MAP} . Since both likelihood and prior are Gaussian, the posterior is Gaussian as well!⁴

$$p(\mathbf{w} \mid \mathcal{D}) = \mathcal{N}(\mathbf{w} \mid \boldsymbol{\mu}, \boldsymbol{\Sigma})$$

where $\boldsymbol{\mu} = \beta \boldsymbol{\Sigma} \boldsymbol{\Phi}^T \mathbf{y}$ and $\boldsymbol{\Sigma}^{-1} = \alpha \mathbf{I} + \beta \boldsymbol{\Phi}^T \boldsymbol{\Phi}$.

Observations

The posterior is Gaussian, so its mode is the mean and $\mathbf{w}_{\text{MAP}} = \boldsymbol{\mu}$

In the limit of an infinitely broad prior $\alpha \rightarrow 0$, $\mathbf{w}_{\text{MAP}} \rightarrow \mathbf{w}_{\text{ML}}$

For $N = 0$, i.e. no data points, the posterior equals the prior

Even though we assume an isotropic prior $p(\mathbf{w})$, the posterior covariance is in general not diagonal

⁴The Gaussian distribution is a *conjugate prior* to itself.

Predictions for new data: MLE and MAP

After observing data $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$, we can compute the MLE/MAP.

Usually, what we are actually interested in is the prediction \hat{y}_{new} for a new data point \mathbf{x}_{new} - the model parameters \mathbf{w} are just a means to achieve this.

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

Plugging in the estimated parameters we get a predictive distribution that lets us make prediction \hat{y}_{new} for new data \mathbf{x}_{new} .

Maximum likelihood given \mathbf{w}_{ML} and β_{ML} :

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

Maximum a posteriori given \mathbf{w}_{MAP} :

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

For MAP we assume β to be known a priori to simplify the calculations.

Predictions for new data: MLE and MAP

After observing data $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$, we can compute the MLE/MAP.

Usually, what we are actually interested in is the prediction \hat{y}_{new} for a new data point \mathbf{x}_{new} - the model parameters \mathbf{w} are just a means to achieve this.

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

Plugging in the estimated parameters we get a **predictive distribution** that lets us make prediction \hat{y}_{new} for new data \mathbf{x}_{new} .

Maximum likelihood given \mathbf{w}_{ML} and β_{ML} :

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

Maximum a posteriori given \mathbf{w}_{MAP} :

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

For MAP we assume β to be known a priori to simplify the calculations.

Posterior predictive distribution

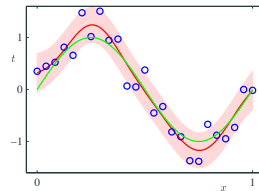
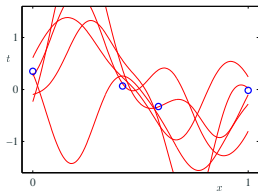
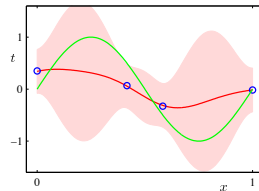
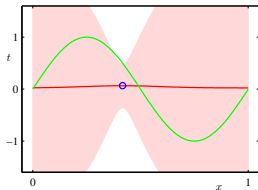
Alternatively, we can use the full posterior distribution $p(\mathbf{w} \mid \mathcal{D})$.

This allows us to compute the **posterior predictive** distribution

$$\begin{aligned} p(\hat{y}_{\text{new}} \mid \mathbf{x}_{\text{new}}, \mathcal{D}) &= \int p(\hat{y}_{\text{new}}, \mathbf{w} \mid \mathbf{x}_{\text{new}}, \mathcal{D}) \, d\mathbf{w} \\ &= \int p(\hat{y}_{\text{new}} \mid \mathbf{x}_{\text{new}}, \mathbf{w}) p(\mathbf{w} \mid \mathcal{D}) \, d\mathbf{w} \\ &= \mathcal{N}(\hat{y}_{\text{new}} \mid \boldsymbol{\mu}^T \phi(\mathbf{x}_{\text{new}}), \beta^{-1} + \phi(\mathbf{x}_{\text{new}})^T \boldsymbol{\Sigma} \phi(\mathbf{x}_{\text{new}})) \end{aligned}$$

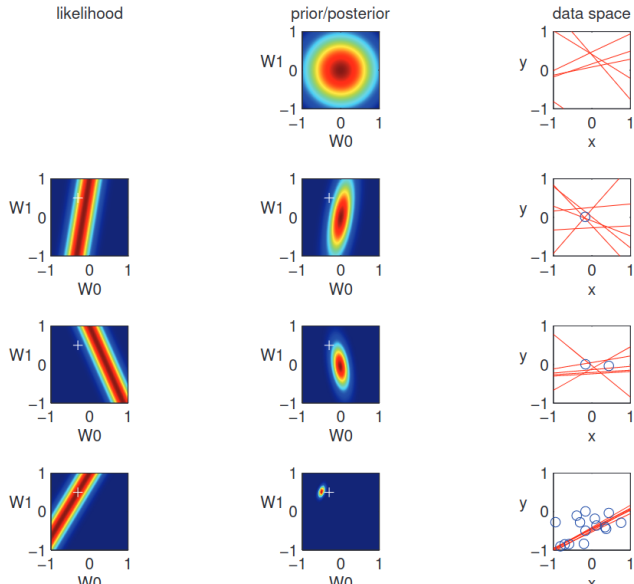
Advantage: We get a data-dependent uncertainty estimate – the variance of the Gaussian now also depends on the input \mathbf{x}_{new} .

Example of posterior predictive distribution



True function (green), observations (blue), mode (dark red), variance (light red).

Sequential Bayesian updating



Optimization-based approaches to regression have probabilistic interpretations:

Least squares regression \iff Maximum likelihood (Slide 27)

Ridge regression \iff Maximum a posteriori (Slide 32)

Even non-linear dependencies can be captured by a model linear w.r.t. weights \boldsymbol{w} by applying a non-linear transformation to the input features (Slide 13).

Penalizing large weights helps to reduce overfitting (Slide 20).

Full Bayesian gives us data-dependent uncertainty estimates (Slide 35).

Main reading

"Pattern Recognition and Machine Learning" by Bishop
[ch. 1.1, 3.1 – 3.6]

Extra reading

"Machine Learning: A Probabilistic Perspective" by Murphy
[ch. 7.2, 7.3, 7.5.1, 7.6.1, 7.6.2]

Slides are based on an older version by G. Jensen, C. Osendorfer, and S. Günnemann. Some figures are from Bishop's "Pattern Recognition and Machine Learning".