



Linear  
Regression

$$x_1, \dots, x_n \in \mathbb{R}^D$$

find a  $f$ , such that  $f(x_i) \approx y_i$

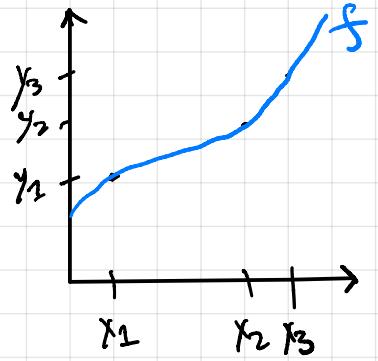
$f$  "generalizes" to new data

$$f(x) = w^\top \Phi(x)$$

↑  
parameters

$\Phi$   $\hat{=}$  basis functions

$$\Phi = \begin{bmatrix} 1 & x_1 & x_1^2 \\ & \vdots & \vdots \\ & x_n & x_n^2 \end{bmatrix}$$



**Problem 1:** Assume that we are given a dataset, where each sample  $x_i$  and regression target  $y_i$  is generated according to the following process

$$x_i \sim \text{Uniform}(-10, 10)$$

$$y_i = ax_i^3 + bx_i^2 + cx_i + d + \epsilon_i, \quad \text{where } \epsilon_i \sim \mathcal{N}(0, 1) \quad \text{and } a, b, c, d \in \mathbb{R}.$$

The 3 regression algorithms below are applied to the given data. Your task is to say what the bias and variance of these models are (low or high). Provide a 1-2 sentence explanation to each of your answers.

- a) Linear regression
- b) Polynomial regression with degree 3
- c) Polynomial regression with degree 10

a) bias: High  
variance: Low

A straight line cannot capture a degree 3 polynomial  
(underfitting) underfitting

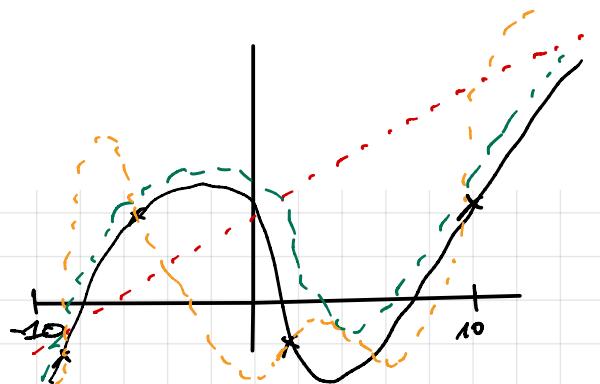
b) bias: Low  
variance: Low

The model is same as the data generating process

We achieve a good fit good fit

c) bias: Low  
variance: High

Since we are using a polynomial regression with a degree much higher compared to the data generating process, the model will overfit the data overfit



**Problem 2:** Given is a training set consisting of samples  $\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N)^T$  with respective regression targets  $\mathbf{y} = (y_1, y_2, \dots, y_N)^T$  where  $\mathbf{x}_i \in \mathbb{R}^D$  and  $y_i \in \mathbb{R}$ .

Alice fits a linear regression model  $f(\mathbf{x}_i) = \mathbf{w}^T \mathbf{x}_i$  to the dataset using the closed form solution for linear regression (normal equations).

Bob has heard that by transforming the inputs  $\mathbf{x}_i$  with a vector-valued function  $\Phi$ , he can fit an alternative function,  $g(\mathbf{x}_i) = \mathbf{v}^T \Phi(\mathbf{x}_i)$ , using the same procedure (solving the normal equations). He decides to use a linear transformation  $\Phi(\mathbf{x}_i) = \mathbf{A}^T \mathbf{x}_i$ , where  $\mathbf{A} \in \mathbb{R}^{D \times D}$  has full rank.

- Show that Bob's procedure will fit the same function as Alice's original procedure, that is  $f(\mathbf{x}) = g(\mathbf{x})$  for all  $\mathbf{x} \in \mathbb{R}^D$  (given that  $\mathbf{w}$  and  $\mathbf{v}$  minimize the training set error).
- Can Bob's procedure lead to a lower training set error than Alice's if the matrix  $\mathbf{A}$  is not invertible? Explain your answer. *No, because it loses dimensions since it hasn't full rank (not invertible)*

a) Alice uses the normal equation directly and obtains:

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

Bob fits the model to the transformed data and obtains:

$$\begin{aligned} \mathbf{v}^* &= [(\mathbf{X} \mathbf{A})^T (\mathbf{X} \mathbf{A})]^{-1} (\mathbf{X} \mathbf{A})^T \mathbf{y} \\ &= (\mathbf{A}^T \mathbf{X}^T \mathbf{X} \mathbf{A})^{-1} \mathbf{A}^T \mathbf{X}^T \mathbf{y} \\ &= \mathbf{A}^{-1} (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{A}^T)^{-1} \mathbf{A}^T \mathbf{X}^T \mathbf{y} \\ &= \mathbf{A}^{-1} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = \mathbf{A}^{-1} \mathbf{w}^* \end{aligned}$$

(Note that  $\Phi$  transforms the column vectors  $\mathbf{x}_i$  via  $\mathbf{A}^T \mathbf{x}_i$  but  $\mathbf{X}$  contains the transposed observations  $\mathbf{x}_i^T$  as rows. Therefore the transformed feature matrix is  $\mathbf{X} \mathbf{A}$ )

$$g(\mathbf{x}_i) = \underbrace{\mathbf{V}^T \Phi(\mathbf{x}_i)}_{\mathbf{A}^T \mathbf{x}_i} = \underbrace{\mathbf{V}^T \mathbf{A}}_{\text{Id}}^{-1} \mathbf{A}^T \mathbf{x}_i = \mathbf{W}^T \mathbf{x}_i = f(\mathbf{x}_i)$$

$$\mathbf{V}^* = \mathbf{A}^{-1} \mathbf{W}^* \quad // \quad \mathbf{V}^{*\top} = (\mathbf{A}^{-1} \mathbf{W}^*)^T = \mathbf{W}^{*\top} \mathbf{A}^{-\top}$$

Any weights  $v^*$  Bob finds are also feasible for Alice by letting  $w = Av^*$ . Therefore Bob can only access a subset of the parameter space and cannot achieve a lower loss value than Alice. (It could still be equal, but it cannot be better). Alice approach is more general

Note that we are only talking about training error in this example, not test error. Bob might manage to find a model that generalizes better than Alice's, but Alice will always be able to fit the training data at least as well as Bob.

a) Alice

$$\frac{1}{2} \sum_{i=1}^N (f(x_i) - y_i)^2 = \frac{1}{2} \sum_{i=1}^N (w^T x_i - y_i)^2$$

Bob loses dimension  
(it does not have full rank)

$$= \frac{1}{2} (x_w - y)^T (x_w - y)$$

$$D_w = \frac{1}{2} (w^T x^T x w - 2 w^T x^T y + y^T y)$$

$$0 \stackrel{!}{=} x^T x_w - x^T y$$

$$\Rightarrow w^* = \underbrace{(x^T x)}^{-1} x^T y$$

$$(AB)^{-1} = B^{-1} A^{-1}$$

Bob

$$\frac{1}{2} \sum_{i=1}^N (w^T \Phi(x_i) - y_i)^2 = \frac{1}{2} \sum_{i=1}^N (w^T A^T x_i - y_i)^2$$

$$= \frac{1}{2} (x_{Aw} - y)^T (x_{Aw} - y) = \frac{1}{2} (w^T A^T x^T x_{Aw} - 2 w^T A^T x^T y + y^T y)$$

$$= A^T x^T x_{Aw} - A^T x^T y \stackrel{!}{=} 0$$

$$A^T x^T x_{Aw} = A x^T y \Rightarrow v^* = (A^T x^T x_{Aw})^{-1} A x^T y \\ = A^{-1} w^*$$

## Least squares regression

**Problem 4:** Let's assume we have a dataset where each datapoint,  $(x_i, y_i)$  is weighted by a scalar factor which we will call  $t_i$ . We will assume that  $t_i > 0$  for all  $i$ . This makes the sum of squares error function look like the following:

$$E_{\text{weighted}}(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N t_i [w^T \phi(x_i) - y_i]^2$$

$$T = \begin{pmatrix} t_1 & t_2 & \dots \end{pmatrix}$$

Find the equation for the value of  $\mathbf{w}$  that minimizes this error function.

Furthermore, explain how this weighting factor,  $t_i$ , can be interpreted in terms of

1) the variance of the noise on the data and

2) data points for which there are exact copies in the dataset.  $\text{yes}$

$$\Sigma_{t_5} = 3$$

$$t_5 \cdot b_5 (3 \cdot b_5)$$

$$\Sigma_{2,3} = \frac{1}{2} [e_2 + e_3 + \underbrace{e_5 + e_5 + e_5}_{t_5 \cdot b_5}]$$

If we define  $T = \text{diag}(t_1, \dots, t_N)$  to be a diagonal matrix with  $t_i$  on the diagonal,

We can write the weighted sum-of-squares cost function in the form

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

Now we follow along the same steps we used to derive the optimal solution for ordinary least squares and arrive at

$$\mathbf{w}^{\text{weighted}} = (\Phi^T T \Phi)^{-1} \Phi^T T \mathbf{y}$$

Can we understand weighted linear regression in a probabilistic context as we did with ordinary least squares?

$$y_i \sim N(w^T \phi(x_i), \beta^{-1})$$

with a common noise precision of  $\beta$ . From this we derived the form of the maximum likelihood error (negative log-likelihood) as

$$E_{\text{ML}}(\mathbf{w}, \beta) = \beta \frac{1}{2} \sum_{i=1}^N (w^T \phi(x_i) - y_i)^2 - \frac{N}{2} \ln \beta + \frac{N}{2} \ln 2 \pi$$

$$\underbrace{\sum_{i=1}^N}_{\mathcal{L}_{\text{LS}}} \quad \underbrace{\ln \beta + \frac{N}{2} \ln 2 \pi}_{\text{const w.r.t } \mathbf{w}}$$

But the least squares error part just stems from the definition of the normal distribution marked with (+) below:

$$N(y_i | w^T \phi(x_i), \beta^{-1}) \propto \exp \left( -\frac{\beta}{2} (w^T \phi(x_i))^2 \right)$$

From this we deduce that weighted least squares is equivalent to probabilistic least squares where we choose  $\beta = t_c$ , in effect modeling  $y_i$  as

$$y_i \sim N(w^T \phi(x_i), t_i^{-1})$$

making the regression targets no longer identically distributed but still dependent. For  $t_i \in \mathbb{N}$ ,  $t_i$  can be regarded as an effective number of replicated observations of data point  $(x_i, y_i)$ .

## Ridge regression

**Problem 5:** Show that ridge regression on a design matrix  $\Phi \in \mathbb{R}^{N \times M}$  with regularization strength  $\lambda$  is equivalent to ordinary least squares regression with an augmented design matrix and target vector

$$\hat{\Phi} = \begin{pmatrix} \Phi \\ \sqrt{\lambda} I_M \end{pmatrix} \quad \text{and} \quad \hat{y} = \begin{pmatrix} y \\ 0_M \end{pmatrix}.$$

Ordinary least squares minimizes  $\frac{1}{2} (\hat{\Phi} w - \hat{y})^\top (\hat{\Phi} w - \hat{y})$ . Plugging in the augmented design matrix and regression target, we get

$$\frac{1}{2} (\hat{\Phi} w - \hat{y})^\top (\hat{\Phi} w - \hat{y}) = \frac{1}{2} (\Phi w - y)^\top (\Phi w - y) +$$

$$\frac{1}{2} (\sqrt{\lambda} I_M w)^\top (\sqrt{\lambda} I_M w)$$

$$= \frac{1}{2} \sum_{i=1}^N (\Phi_i w - y_i)^2 + \frac{\lambda}{2} \|w\|_2^2$$

which is equal to the ridge regression loss function

## Implementation

**Problem 6:** John Doe is a data scientist, and he wants to fit a polynomial regression model to his data. For this, he needs to choose the degree of the polynomial that works best for his problem.

Unfortunately, John hasn't attended IN2064, so he writes the following code for choosing the optimal degree of the polynomial:

```
X, y = load_data()
best_error = -1
best_degree = None

for degree in range(1, 50):
    w = fit_polynomial_regression(X, y, degree)
    y_predicted = predict_polynomial_regression(X, w, degree)
    error = compute_mean_squared_error(y, y_predicted)
    if (error <= best_error) or (best_error == -1):
        best_error = error
        best_degree = degree

print("Best degree is " + str(best_degree))
```

Assume that the functions are implemented correctly and do what their name suggests.

- Explain briefly why this code doesn't do what it's supposed to do.

Output: Best degree is 49

Error on training set always goes down when we use higher degree polynomial (unless it's already 0, then it stays at 0) overfitting

- Describe a possible way to fix the problem with this code. (You don't need to write any code, just describe the approach.)

Split data into train and validation sets. Choose the degree that achieves the lowest mean squared error on the validation set (not on the training set)

It does not work well on unseen datapoints

## Bayesian linear regression

Bishop 3.12

In the lecture we made the assumption that we already knew the precision (inverse variance) for our Gaussian distributions. What about when we don't know the precision and we need to put a prior on that as well as our Gaussian prior that we already have on the weights of the model?

**Problem 7:** It turns out that the conjugate prior for the situation when we have an unknown mean and unknown precision is a normal-gamma distribution (See section 2.3.6 in Bishop). This means that if our likelihood is as follows:

$$p(\mathbf{y} | \Phi, \mathbf{w}, \beta) = \prod_{i=1}^N \mathcal{N}(y_i | \mathbf{w}^T \phi(\mathbf{x}_i), \beta^{-1})$$

Then the conjugate prior for both  $\mathbf{w}$  and  $\beta$  is

$$p(\mathbf{w}, \beta) = \mathcal{N}(\mathbf{w} | \mathbf{m}_0, \beta^{-1} \mathbf{S}_0) \text{Gamma}(\beta | a_0, b_0)$$

Show that the posterior distribution takes the same form as the prior, i.e.

$$p(\mathbf{w}, \beta | \mathcal{D}) = \mathcal{N}(\mathbf{w} | \mathbf{m}_N, \beta^{-1} \mathbf{S}_N) \text{Gamma}(\beta | a_N, b_N)$$

Also be sure to give the expressions for  $\mathbf{m}_N$ ,  $\mathbf{S}_N$ ,  $a_N$ , and  $b_N$ .

*Hint:* Expand  $\log p(\mathbf{w}, \beta | \mathcal{D})$  once with the prior and likelihood and once with the presumed posterior form. The resulting expressions have to be equal, so you should be able to match all terms in the two expansions against each other and then read off the parameters of the posterior distribution.

It is easiest to work in log space. The log of the posterior distribution is given by

$$p(\mathbf{w}, \beta | \mathcal{D}) = \frac{p(\mathcal{D} | \mathbf{w}, \beta) \cdot p(\mathbf{w}, \beta)}{p(\mathcal{D})}$$

$$\begin{aligned} \log p(\mathbf{w}, \beta | \mathcal{D}) &= \log p(\mathbf{w}, \beta) + \sum_{i=1}^N \log p(y_i | \mathbf{w}^T \phi(\mathbf{x}_i), \beta^{-1}) + \text{const} \\ &= \frac{M}{2} \log \beta - \frac{\beta}{2} (\mathbf{w} - \mathbf{m}_0)^T \mathbf{S}_0^{-1} (\mathbf{w} - \mathbf{m}_0) - \frac{1}{\beta} \log \beta + (a_0 - 1) \log \beta + \\ &\quad \frac{N}{2} \log \beta - \frac{\beta}{2} \sum_{i=1}^N \{ \mathbf{w}^T \phi(\mathbf{x}_i) - y_i \}^2 + \text{const}. \end{aligned}$$

$$\log p(\mathbf{w}, \beta | \mathcal{D}) = \log p(\mathbf{w}, \beta | \mathcal{D}) + \log p(\beta | \mathcal{D})$$

The general form of  $\log p(w | \beta, D)$  is given by

$$\frac{\beta}{2} (w - m_N)^T S_N^{-1} (w - m_N) + \frac{M}{2} \log \beta + \text{const}$$

which we expand to

$$-\frac{\beta}{2} w^T S_N^{-1} w + \beta m_N^T S_N^{-1} w - \frac{\beta}{2} m_N^T S_N^{-1} m_N + \frac{M}{2} \log \beta + \text{const}$$

**Problem 8:** Derive the closed form solution for ridge regression error function

$$E_{\text{ridge}}(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N (\mathbf{w}^T \phi(\mathbf{x}_i) - y_i)^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}$$

Additionally, discuss the scenario when the number of training samples  $N$  is smaller than the number of basis functions  $M$ . What computational issues arise in this case? How does regularization address them?

$$\begin{aligned} E_{\text{ridge}}(\mathbf{w}) &= \frac{1}{2} \sum_{i=1}^N (\mathbf{w}^T \phi(\mathbf{x}_i) - y_i)^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w} \\ &= \frac{1}{2} (\Phi \mathbf{w} - \mathbf{y})^T (\Phi \mathbf{w} - \mathbf{y}) + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w} \end{aligned}$$

Taking the gradient

$$\begin{aligned} \nabla_{\mathbf{w}} E_{\text{ridge}}(\mathbf{w}) &= \Phi^T \Phi \mathbf{w} - \Phi^T \mathbf{y} + \lambda \mathbf{w} \\ &= (\Phi^T \Phi + \lambda I) \mathbf{w} - \Phi^T \mathbf{y} \end{aligned}$$

Set it to zero:

$$\begin{aligned} (\Phi^T \Phi + \lambda I) \mathbf{w} &= \Phi^T \mathbf{y} \\ \mathbf{w} &= (\underbrace{\Phi^T \Phi + \lambda I}_{\text{if this matrix is not invertible}})^{-1} \Phi^T \mathbf{y} \end{aligned}$$

If  $N < M$ , the covariance matrix  $\Phi^T \Phi \in \mathbb{R}^{N \times N}$  will be singular, therefore not invertible. (this may happen even if  $N \geq M$ , e.g. when some features are correlated).

When regularization is used,  $\lambda I$  is added to the covariance matrix, thus fixing the potential degeneracy issue and making the problem tractable.

## Comparison of Linear Regression Models *Exam exercise*

**Problem 9:** We want to perform regression on a dataset consisting of  $N$  samples  $\mathbf{x}_i \in \mathbb{R}^D$  with corresponding targets  $y_i \in \mathbb{R}$  (represented compactly as  $\mathbf{X} \in \mathbb{R}^{N \times D}$  and  $\mathbf{y} \in \mathbb{R}^N$ ).

Assume that we have fitted an  $L_2$ -regularized linear regression model and obtained the optimal weight vector  $\mathbf{w}^* \in \mathbb{R}^D$  as

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

Note that there is no bias term.

Now, assume that we obtained a new data matrix  $\mathbf{X}_{new}$  by scaling all samples by the same positive factor  $a \in (0, \infty)$ . That is,  $\mathbf{X}_{new} = a\mathbf{X}$  (and respectively  $\mathbf{x}_i^{new} = a\mathbf{x}_i$ ).

- a) Find the weight vector  $\mathbf{w}_{new}$  that will produce the same predictions on  $\mathbf{X}_{new}$  as  $\mathbf{w}^*$  produces on  $\mathbf{x}$

a) Predictions of a linear regression model are generated as

$$\hat{\mathbf{y}} = \mathbf{w}^T \mathbf{x}.$$

$$\mathbf{w}^T \mathbf{x}_i = \mathbf{x}_{new}^T \mathbf{x}_i^{new} \text{ or equivalently } \mathbf{w}^T \mathbf{x}_i = \mathbf{w}_{new}^T a \mathbf{x}_i$$

$$\text{Solving for } \mathbf{w}_{new} \text{ we get } \mathbf{w}_{new} = \mathbf{w}^*/a$$

- b) Find the regularization factor  $\lambda_{new} \in \mathbb{R}$ , such that the solution  $w_{new}^*$  of the new  $L_2$ -regularized linear regression problem

$$w_{new}^* = \arg \min_w \frac{1}{2} \sum_{i=1}^N (w^T X_i^{new} - y_i)^2 + \frac{\lambda_{new}}{2} w^T w$$

will produce the same predictions on  $X_{new}$  as  $w^*$  produces on  $X$ .

Provide a mathematical justification for your answer.

The closed form solution for  $w^*$  on the original data  $X$  is

$$w^* = (X^T X + \lambda I)^{-1} X^T y$$

The closed form solution for  $w_{new}^*$  on the new data  $X_{new}$  is

$$\begin{aligned} w_{new}^* &= (X_{new}^T X_{new} + \lambda_{new} I)^{-1} X_{new}^T y \\ &= \alpha (a^2 X^T X + \lambda_{new} I)^{-1} X^T y \end{aligned}$$

by setting  $\lambda_{new} = a^2 \lambda$ , we get

$$\begin{aligned} &= \alpha (a^2 X^T X + a^2 \lambda I)^{-1} X^T y \\ &= \frac{1}{a} (X^T X + \lambda I)^{-1} X^T y \\ &= \frac{1}{a} w^* \end{aligned}$$

which (according to our answer in part (a)) will produce the same predictions on  $X_{new}$  as  $w^*$  does on  $X$ , as desired

Equivalent solution

$$\begin{aligned} w_{new}^* &\stackrel{!}{=} \frac{w^*}{a} = \frac{1}{a} \arg \min_w \frac{1}{2} \sum_{i=1}^N (w^T x_i - y_i)^2 + \frac{\lambda}{2} w^T w \\ &= \frac{1}{a} \arg \min_w \frac{1}{2} \sum_{i=1}^N \left( \frac{w^T}{a} a x_i - y_i \right)^2 + \frac{a^2 \lambda}{2} \frac{w^T}{a} \frac{w}{a} \\ &= \frac{a}{a} \arg \min_{\substack{w \\ w_{new} = \frac{w}{a}}} \frac{1}{2} \sum_{i=1}^N (w_{new}^T X_i^{new} - y_i)^2 + \frac{a^2 \lambda}{2} w_{new}^T w_{new} \\ &\stackrel{!}{=} w_{new}^* = \arg \min_{w_{new}} \frac{1}{2} \sum_{i=1}^N (w_{new}^T X_i^{new} - y_i)^2 + \frac{\lambda_{new}}{2} w_{new}^T w_{new} \end{aligned}$$

We have to  
set  $\lambda_{new} = a^2 \lambda$