

# Machine Learning

## Linear Regression

Prof. Matthias Hein

Machine Learning Group  
Department of Mathematics and Computer Science  
Saarland University, Saarbrücken, Germany

**Lecture 6, 8.11.2013**

**Linear methods:** most simple regression and classification techniques.

- easy interpretation: feature has a high influence if it has a large weight.
- linear methods: have possibly high bias but low variance  $\Rightarrow$  can be fit already with only a few training points.
- often competitive with non-linear methods in high dimensions,
- Using transformations of the input features (**basis functions**) one can easily generate non-linear functions in the input space.

**Important:** Linear methods are *linear* in the parameters, but not necessarily linear in the original input features.

## Risk of squared loss:

$$\mathbb{E}[(Y - f(X))^2] = \mathbb{E}[\mathbb{E}[(Y - f(X))^2 | X]].$$

## Bayes optimal function:

$$f(x) = \mathbb{E}[Y | X = x].$$

## Definition

Given a training sample  $T_n = (X_i, Y_i)_{i=1}^n$  with  $X_i \in \mathcal{X}$  and  $Y_i \in \mathbb{R}$  and a function space  $\mathcal{F}$  we define **least squares regression** as the mapping  $\mathcal{A}: T_n \rightarrow \mathcal{F}$  with,

$$T_n \mapsto f_n = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (Y_i - f(X_i))^2.$$

**Linear least squares regression:** used **Function class**:

$$\mathcal{F} = \left\{ f \mid f(x) = \sum_{i=1}^d w_i x_i + b = \langle w, x \rangle + b, \quad w \in \mathbb{R}^d, b \in \mathbb{R} \right\}.$$

**Notation:**

- $w$  is the **weight vector**,
- summarize the outputs  $(Y_i)_{i=1}^n$  into a column vector  $Y \in \mathbb{R}^n$  and the inputs vectors  $(X_i)_{i=1}^n$  into a matrix  $X \in \mathbb{R}^{n \times d}$ ,

$$Y = \begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix}, \quad X = \begin{pmatrix} X_{11} & \dots & X_{1d} \\ \vdots & & \vdots \\ X_{n1} & \dots & X_{nd} \end{pmatrix},$$

where  $X$  is called the **design matrix**.

- **Convention in the lecture:** vectors are always column vectors !

# Least squares regression III

## Constant term in the function class $\mathcal{F}$ :

- Add extra dimension to input vector to integrate constant term in the function class,

$$X'_i = (X_{i1}, \dots, X_{id}, 1) \quad \text{or} \quad X'_{i(d+1)} = 1, \quad \forall i.$$

An affine function is characterized by the weight vector  $w$ ,

$$w \in \mathbb{R}^{d+1}, \quad f(X'_i) = \langle w, X'_i \rangle = \sum_{j=1}^{d+1} w_j X'_{ij} = \sum_{j=1}^d w_j X_{ij} + w_{d+1}.$$

## Linear least squares regression:

$$w_n = \arg \min_{w \in \mathbb{R}^{d+1}} \frac{1}{n} \sum_{i=1}^n (Y_i - \langle X_i, w \rangle)^2$$

**Convention: we make the constant  $b$  explicit in the lecture**

## Proposition

Let  $X \in \mathbb{R}^{n \times d}$ . The solution  $w_n$  of linear least squares regression is given by

$$w_n = (X^T X)^{-1} X^T Y,$$

where the inverse  $(X^T X)^{-1}$  exists if  $X$  has rank  $d$ . If  $X$  has not rank  $d$ , then  $(X^T X)^{-1}$  has to be understood in the sense of a generalized inverse. In this case the solution is not unique but if  $w_n^1, w_n^2$  are two solutions, then the predictions agree on the training data

$$f_{w_n^1}(X_i) = \langle w_n^1, X_i \rangle = \langle w_n^2, X_i \rangle = f_{w_n^2}(X_i), \quad \text{for all } i = 1, \dots, n.$$

# Least squares regression V

**Proof:** Objective function of the optimization problem with  $w \in \mathbb{R}^d$ ,

$$O_{LLSR}(w) = \frac{1}{n} \sum_{i=1}^n (Y_i - \langle w, X_i \rangle)^2 = \frac{1}{n} \|Y - Xw\|^2.$$

Taking the derivative with respect to  $w$ ,

$$\nabla_w O_{LLSR} = -\frac{2}{n} X^T (Y - Xw).$$

The necessary condition for an extremum of  $O_{LLSR}$  is therefore

$$\frac{2}{n} X^T (Y - Xw) = 0 \quad \implies \quad X^T Y = (X^T X)w \quad w_n = (X^T X)^{-1} X^T Y$$

Hessian of the objective function  $\frac{2}{n} X^T X \Rightarrow$  positive-definite if  $X$  has rank  $d$ . If  $X$  has rank smaller than  $d$ , then  $w_n$  defined using the generalized inverse is a solution and every  $w = w_n + v$  where  $v$  is orthogonal to the subspace  $\text{Span}\{X_1, \dots, X_n\}$  is another solution.

# The pseudo-inverse

$(X^T X)^{-1} X^T$  is the Moore-Penrose **pseudo inverse** of  $X$  if  $X$  has rank  $d$ .

## Definition

Let  $A \in \mathbb{R}^{m \times n}$  with rank  $r \leq \min\{m, n\}$ . Then the **pseudo-inverse**  $A^+$  of  $A$  is defined as

$$A^+ = \arg \min_{B \in \mathbb{R}^{n \times m}} \|AB - \mathbb{1}_m\|_F^2,$$

where  $\|\cdot\|_F$  is the **Frobenius norm** ( $\|A\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n A_{ij}^2}$ ) and  $\mathbb{1}_m$  the identity matrix in  $\mathbb{R}^m$ .

Let  $A$  be a square matrix which is invertible, then

$$(A^T A)^{-1} A^T = A^{-1} (A^T)^{-1} A^T = A^{-1}.$$



# The pseudo-inverse II

The **singular value decomposition** of  $A \in \mathbb{R}^{m \times n}$ ,

$$A = U \Sigma V^T,$$

- $U$  is an orthogonal matrix  $U = (u_1, \dots, u_m) \in \mathbb{R}^{m \times m}$ , that is  $U^T U = \mathbb{1}_m$ ,
- $V$  is an orthogonal matrix  $V = (v_1, \dots, v_n) \in \mathbb{R}^{n \times n}$ , that is  $V^T V = \mathbb{1}_n$ ,
- $\Sigma \in \mathbb{R}^{m \times n}$  with  $\Sigma_{ij} = \begin{cases} \sigma_i & \text{if } i = j \text{ and } i \leq r, \\ 0 & \text{otherwise} \end{cases}$ .

The  $\sigma_i > 0$ ,  $i = 1, \dots, r$  are the **singular values** of  $A$ .

The **pseudo inverse**  $A^+$  is then given by

$$A^+ = V \Sigma^+ U^T,$$

where  $\Sigma^+ \in \mathbb{R}^{n \times m}$  is defined as  $\Sigma_{ij}^+ = \begin{cases} 1/\sigma_i & \text{if } i = j \text{ and } i \leq r, \\ 0 & \text{otherwise} \end{cases}$ .

The **pseudo inverse**  $A^+$  is then given by

$$A^+ = V\Sigma^+U^T,$$

where  $\Sigma^+ \in \mathbb{R}^{n \times m}$  is given by  $\Sigma_{ij}^+ = \begin{cases} 1/\sigma_i & \text{if } i = j \text{ and } i \leq r, \\ 0 & \text{otherwise} \end{cases}$ .

Let  $A \in \mathbb{R}^{n \times m}$ . Given that  $m \leq n$  and  $\text{ran}(A) = m$ , one can write the pseudo inverse  $A^+$  as  $A^+ = (A^T A)^{-1} A^T$ ,

$$\begin{aligned} (A^T A)^{-1} A^T &= (V\Sigma^T U^T U \Sigma V^T)^{-1} V\Sigma^T U^T = (V\Sigma^T \Sigma V^T)^{-1} V\Sigma^T U^T \\ &= V(\Sigma^T \Sigma)^{-1} V^T V \Sigma^T U^T = V(\Sigma^T \Sigma)^{-1} \Sigma^T U^T = V\Sigma^+ U^T. \end{aligned}$$

## Basis functions/Feature maps:

- map the input  $x \rightarrow \phi(x)$ ,  
 $\mathcal{X} = \mathbb{R} : x, x^2, x^3, \dots$  (polynomials),  
 $\mathcal{X} = [0, 2\pi] : \sin(x), \cos(x), \sin(2x), \cos(2x), \dots$  (Fourier basis).

Fixed, pre-defined set of  $D$  **basis functions**,  $\phi_i : \mathbb{R}^d \rightarrow \mathbb{R}$ , we define the function space

$$\mathcal{F} = \left\{ f : \mathbb{R}^d \rightarrow \mathbb{R}, f(x) = \sum_{i=1}^D w_i \phi_i(x) \mid w \in \mathbb{R}^D \right\}.$$

**Advantage:** explicit integration of prior knowledge possible.

**Generalized design matrix:**  $\Phi \in \mathbb{R}^{n \times D}$ ,

$$\Phi = \begin{pmatrix} \phi_1(X_1) & \dots & \phi_D(X_1) \\ \vdots & & \vdots \\ \phi_1(X_n) & \dots & \phi_D(X_n) \end{pmatrix},$$

**Least squares regression problem:**

$$w_n = \arg \min_{w \in \mathbb{R}^D} \frac{1}{n} \sum_{i=1}^n (Y_i - \langle w, \phi(X_i) \rangle)^2 = \frac{1}{n} \|Y - \Phi w\|^2,$$

with solution

$$w_n = (\Phi^T \Phi)^{-1} \Phi^T Y,$$

where the matrix  $(\Phi^T \Phi)^{-1} \Phi^T \in \mathbb{R}^{D \times n}$  is the pseudo-inverse of  $\Phi$ .

## Properties:

- The final function,  $f(x) = \langle w_n, \phi(x) \rangle = \sum_{i=1}^D w_i \phi_i(x)$ , is linear in the parameter  $w$ ,
- allows direct modeling of prior knowledge,
- function space  $\mathcal{F} = \left\{ f(x) = \sum_{i=1}^D w_i \phi_i(x) \mid w \in \mathbb{R}^D \right\}$  is  $D$ -dimensional,
- Problem: want to model all polynomials in  $\mathbb{R}^d$ ,  $d$  polynomials of degree one (linear functions),  $\frac{d(d+1)}{2}$  polynomials of degree two, .... Set of basis functions increases rapidly with the dimension  $d \Rightarrow$  not practical.

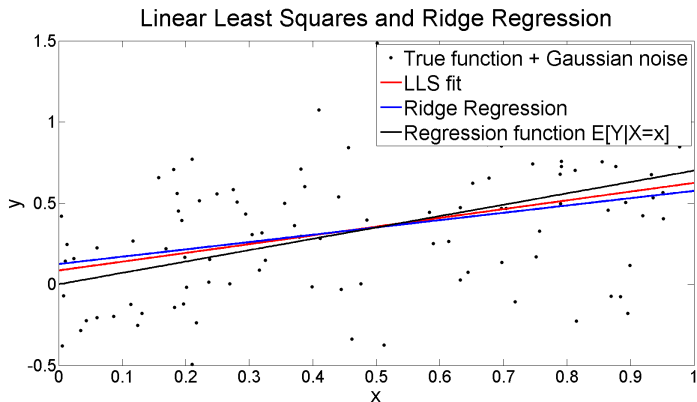
## Ridge regression:

- **Motivation:** originally: add small ridge to the solution so that it becomes unique,  
today: regularized version of the least squares problem.
- **Function space:**  $\mathcal{F} = \left\{ f(x) = \sum_{i=1}^D w_i \phi_i(x) \mid w \in \mathbb{R}^D \right\}$
- **Loss:** squared loss
- **Regularizer:**  $\Omega(w) = \sum_{i=1}^D w_i^2 = \|w\|_2^2$ .

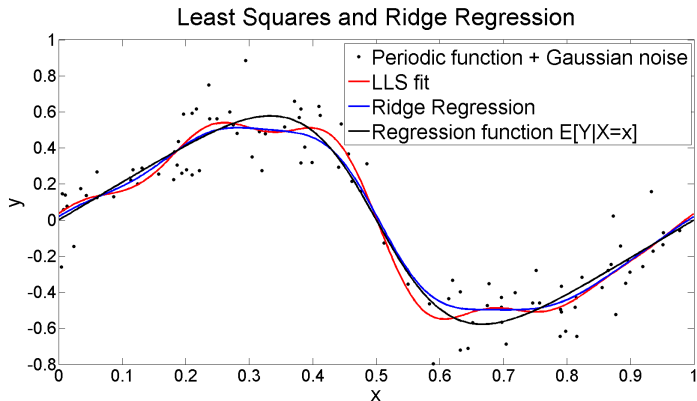
## Definition

Given sample  $T_n = (X_i, Y_i)_{i=1}^n$ , **ridge regression** is defined as the mapping  $\mathcal{A} : T_n \rightarrow \mathcal{F}$  with,

$$T_n \mapsto f_n = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (Y_i - \langle w, \phi_i(x) \rangle)^2 + \lambda \sum_{i=1}^D w_i^2.$$



**Figure :** Linear least squares regression versus linear ridge regression. The regression function is linear.



**Figure :** Comparison of least squares and ridge regression using a set of periodic basis functions.



## Solution of ridge regression:

$$w_{n,\lambda} = (\Phi^T \Phi + \lambda \mathbb{1}_D)^{-1} \Phi^T Y.$$

## Properties:

- solution  $w_{n,\lambda}$  exists and is unique,
- regularizer  $\Omega(w) = \|w\|^2$  corresponds to

$$p(w) \propto e^{-\Omega(w)} = e^{-\|w\|^2}.$$

as a prior for maximum a posteriori (MAP) estimation

# Geometric interpretation

**Linear least squares regression:** use SVD of  $X$ ,  $X = U\Sigma V^T$ , where  $\text{rank } \Sigma = r$ ,

$$Xw_n = X(X^T X)^{-1} X^T Y = U\Sigma V^T V^T (\Sigma^+)^2 V^T V \Sigma^T U^T Y = \sum_{i=1}^r u_i \langle u_i, Y \rangle.$$

**Ridge regression:**

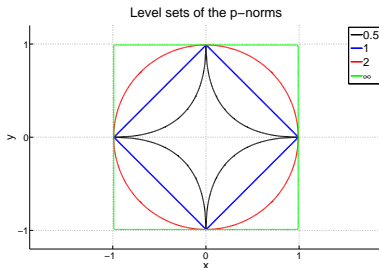
$$Xw_{n,\lambda} = X(X^T X + \lambda \mathbb{1}_d)^{-1} X^T Y = U F(\Sigma) U^T Y = \sum_{i=1}^r u_i \frac{\sigma_i^2}{\sigma_i^2 + \lambda} \langle u_i, Y \rangle,$$

where  $F(\Sigma) = \begin{cases} \frac{\sigma_i^2}{\sigma_i^2 + \lambda} & \text{if } i = j \text{ and } i \leq r, \\ 0 & \text{otherwise} \end{cases}$ ,  $\sigma_i$  are singular values of  $X$ .

- outputs are projected on the basis spanned by  $U$ ,
- The directions  $u_i = \frac{1}{\sigma_i} X v_i$  correspond to the (mapped) eigenvectors  $v_i$  of the covariance matrix  $C_{ij} = X^T X$  if  $X$  is centered.

# The lasso - Least Squares with $L_1$ -Regularization I

**Other regularization functionals:**  $\Omega(w) = \sum_{i=1}^n |w_i|^p = \|w\|_p^p$ .  
 $\Rightarrow$   $L_2$ -norm is the only **isotropic** norm in the family of  $p$ -norms !



**Figure :** The level set  $\|w\|_p = 1$  of the  $p$ -norms. Note that the  $\|\cdot\|_p$  is only a norm for  $p \geq 1$ , in which case the unit-ball is a convex set. Clearly for  $p = 0.5$  the “unit-ball” is not convex.

## Definition

Given a training sample  $T_n = (X_i, Y_i)_{i=1}^n$  with  $X_i \in \mathcal{X}$  and  $Y_i \in \mathbb{R}$  and the function space  $\mathcal{F} = \{\sum_{j=1}^D w_j \phi_j(x) \mid w \in \mathbb{R}^D\}$  we define **the lasso** as the mapping  $\mathcal{A} : T_n \rightarrow \mathcal{F}$  with,

$$T_n \mapsto w_n = \arg \min_{w \in \mathbb{R}^D} \frac{1}{n} \sum_{i=1}^n (Y_i - \langle w, \Phi(X_i) \rangle)^2 + \lambda \sum_{i=1}^D |w_i|.$$

## Motivation:

- $L_1$ -norm induces **sparsity** (a lot of components  $w_i$  of  $w$  are zero).  
**Why ?** The “zero norm” (not really a norm) enforces directly sparsity:

$$\|w\|_0 = \sum_{i=1}^D \mathbb{1}_{w_i \neq 0}.$$

$L_1$ -norm is the norm which is “closest” to the “zero norm” !

**Sparsity is good :** less storage, faster evaluation  $f(x) = \langle w, x \rangle$ ,  
feature selection

## Motivation:

- $L_1$ -norm induces **sparsity** (a lot of components  $w_i$  of  $w$  are zero).  
**Why ?** The “zero norm” (not really a norm) enforces directly sparsity:

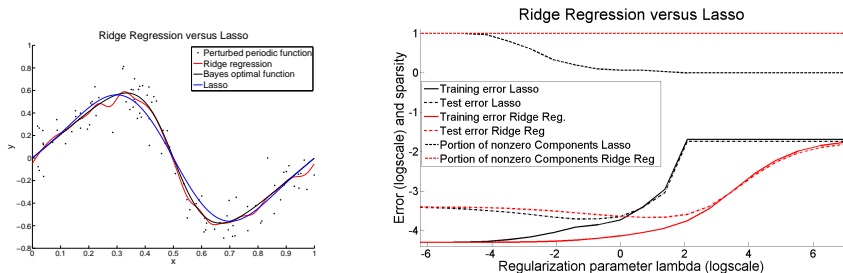
$$\|w\|_0 = \sum_{i=1}^D \mathbb{1}_{w_i \neq 0}.$$

$L_1$ -norm is the norm which is “closest” to the “zero norm” !

**Sparsity is good** : less storage, faster evaluation  $f(x) = \langle w, x \rangle$ ,  
feature selection

- $\|w\|_2^2$  penalizes large weights heavily  $\Rightarrow$  preference for small weights in all directions. (regularizer is **isotropic**)  
 $\|w\|_1$  penalizes large and small weights “equally”  $\Rightarrow$  produces often large weights in few directions.

# Comparison: lasso and ridge regression



**Figure :** **Left:** Perturbed training data and regression function in black, we show the solution of ridge regression in blue and of Lasso in red for  $\lambda = 1$ , **Right:** Behavior of training and test error and number of non-zero components of the weight vector as a function of the regularization parameter  $\lambda$ .

# Bias and variance of estimators

Solutions  $w_n$  of least squares or ridge regression are estimators for the optimal parameter  $w^*$  (Bayes optimal **linear** function for the squared loss),

$$w^* = \arg \min_{w \in \mathbb{R}^d} \mathbb{E}[(Y - \langle w, X \rangle)^2] = \mathbb{E}\left[\left(Y - \sum_{i=1}^d w_i X_i\right)^2\right],$$

The solution can be derived as ( $X$  is a row vector !):

$$w^* = \left(\mathbb{E}[X^T X]\right)^{-1} \mathbb{E}[X^T Y].$$

The empirical solutions  $w_n$  depend on the training sample  $T = (X_i, Y_i)$ .

## Questions:

- Is the average estimator  $w_n$  over training samples of size  $n$  equal to the optimal  $w^*$  ?
- How much does the estimator  $w_n$  fluctuate around its average value over all possible training samples from  $P$  of size  $n$  ?



## Definition

Given a sample  $T = (X_i)_{i=1}^n$  and an estimate (also called statistics)  $f_n : T \rightarrow \mathbb{R}$  of a quantity  $f \in \mathbb{R}$  the **bias** of  $f_n$  is defined as

$$\text{Bias } f_n = \mathbb{E}_T[f_n] - f,$$

the difference of the expectation of  $f_n$  over all training sets  $T$  (all possible i.i.d. training sets of size  $n$ ) and the true quantity  $f$ .

- The estimator  $f_n$  is said to be **unbiased** if the bias is zero.
- It is **asymptotically unbiased** if  $\lim_{n \rightarrow \infty} \text{Bias } f_n = 0$ .

The **variance** of  $f_n$  is defined as,

$$\text{Var } f_n = \mathbb{E}_T[(f_n - \mathbb{E}_T[f_n])^2].$$

## Examples for bias and variance:

- **The empirical mean**  $\mathbb{E}_{P_n}[X] = \frac{1}{n} \sum_{i=1}^n X_i$  is an **estimator of the true mean**  $\mathbb{E}[X] = \mathbb{E}_P[X]$ .

$$\mathbb{E}_T \left[ \frac{1}{n} \sum_{i=1}^n X_i \right] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{X_i} [X_i] = \frac{1}{n} n \mathbb{E}[X] = \mathbb{E}[X] \implies \text{unbiased!}$$

- **empirical variance**  $\text{Var}_{P_n}[X] = \frac{1}{n} \sum_{i=1}^n (X_i - \mathbb{E}_{P_n}[X])^2$  as an estimator of the **true variance**  $\text{Var}_P[X] = \text{Var}[X]$ .

$$\mathbb{E}_T [\text{Var}_{P_n}[X]] = \frac{n-1}{n} \text{Var}[X] \implies \text{biased! underestimation!}$$

The estimator  $\frac{1}{n-1} \sum_{i=1}^n (X_i - \mathbb{E}_{P_n}[X])^2$  for the variance of  $X$  is unbiased.

The risk  $R(f_n)$ , the expected squared loss, of the estimator  $f_n$ :

$$\begin{aligned} R(f_n) &= \mathbb{E}[(Y - f_n(X))^2] = \mathbb{E}[\mathbb{E}[(Y - f_n(X))^2|X]] \\ &= \mathbb{E}[\mathbb{E}[(Y - \mathbb{E}[Y|X] + \mathbb{E}[Y|X] - f_n(X))^2|X]] \\ &= \mathbb{E}[\mathbb{E}[(Y - \mathbb{E}[Y|X])^2|X]] + \mathbb{E}[\mathbb{E}[(\mathbb{E}[Y|X] - f_n(X))^2|X]] \\ &\quad + 2\mathbb{E}[\mathbb{E}[(Y - \mathbb{E}[Y|X])(\mathbb{E}[Y|X] - f_n(X))|X]], \\ &= \mathbb{E}[(Y - \mathbb{E}[Y|X])^2] + \mathbb{E}[(\mathbb{E}[Y|X] - f_n(X))^2] \end{aligned}$$

## Interpretation:

- The first term is the **Bayes optimal risk** (often also called noise term),  
 $\eta(x) = \mathbb{E}[Y|X = x]$  is the Bayes optimal function for the squared loss.
- The second term measures the **deviation of  $f_n$  from the Bayes optimal function**. It is a random quantity since  $f_n$  depends on the training data !

**Expected risk**  $\mathbb{E}_T[R(f_n)]$  over all possible training sets  $T$ :

$$\mathbb{E}_T[R(f_n)] = \mathbb{E}[(Y - \eta(X))^2] + \mathbb{E}_T[\mathbb{E}_X[(\eta(X) - f_n(X))^2]],$$

The first term is constant !

$$\begin{aligned}\mathbb{E}_T[(f_n(x) - \eta(x))^2] &= \mathbb{E}_T[(f_n(x) - \mathbb{E}_T f_n(x) + \mathbb{E}_T f_n(x) - \eta(x))^2] \\&= \mathbb{E}_T[(f_n(x) - \mathbb{E}_T[f_n(x)])^2] + \mathbb{E}_T[(\mathbb{E}_T[f_n(x)] - \eta(x))^2] \\&\quad + 2\mathbb{E}_T[(f_n(x) - \mathbb{E}_T f_n(x))(\mathbb{E}_T f_n(x) - \eta(x))] \\&= \mathbb{E}_T[(f_n(x) - \mathbb{E}_T[f_n(x)])^2] + (\mathbb{E}_T[f_n(x)] - \eta(x))^2 \\&= \text{Var } f_n(x) + (\text{Bias } f_n(x))^2,\end{aligned}$$

## (Noise)-Bias-Variance-Decomposition:

$$\mathbb{E}_T[R(f_n)] = \mathbb{E}[(Y - \eta(X))^2] + \mathbb{E}[(\text{Bias } f_n(X))^2] + \mathbb{E}[\text{Var } f_n(X)],$$

where

- **Noise term at  $x$ :**  $\mathbb{E}[(Y - \eta(X))^2 | X = x],$
- **Variance of  $f_n$ :**  $\text{Var } f_n(x) = \mathbb{E}_T[(f_n(x) - \mathbb{E}_T[f_n(x)])^2],$
- **Bias of  $f_n$ :**  $\text{Bias } f_n(x) = \mathbb{E}_T[f_n(x)] - \eta(x),$

## (Noise)-Bias-Variance-Decomposition:

$$\mathbb{E}_T[R(f_n)] = \mathbb{E}[(Y - \eta(X))^2] + \mathbb{E}[(\text{Bias } f_n(X))^2] + \mathbb{E}[\text{Var } f_n(X)],$$

where

- **Noise term at  $x$ :**  $\mathbb{E}[(Y - \eta(X))^2 | X = x],$
- **Variance of  $f_n$ :**  $\text{Var } f_n(x) = \mathbb{E}_T[(f_n(x) - \mathbb{E}_T[f_n(x)])^2],$
- **Bias of  $f_n$ :**  $\text{Bias } f_n(x) = \mathbb{E}_T[f_n(x)] - \eta(x),$

expected loss = noise + variance + squared bias.

Trade-off between **bias** and **variance**  
corresponds to  
Trade-off between **overfitting** and **underfitting**.

## Bias-Variance-Decomposition for the Least-Squares estimator:

$f_n = \langle w_n, x \rangle \Rightarrow$  express bias and variance of  $f_n$  via the bias and covariance of  $w_n$ ,

$$\begin{aligned}\text{Bias } f_n(x) &= \mathbb{E}_T[f_n(x)] - f^*(x) = \mathbb{E}_T[\langle w_n, x \rangle] - \langle w^*, x \rangle \\ &= \langle \mathbb{E}_T[w_n] - w^*, x \rangle \\ &= \langle \text{Bias } w_n, x \rangle,\end{aligned}$$

$$\begin{aligned}\text{Var } f_n(x) &= \mathbb{E}_T[(f_n(x) - \mathbb{E}_T[f_n(x)])^2] = \mathbb{E}_T[(\langle w_n, x \rangle - \langle \mathbb{E}_T[w_n], x \rangle)^2] \\ &= \mathbb{E}_T[\langle w_n - \mathbb{E}_T[w_n], x \rangle^2] \\ &= \mathbb{E}_T\left[\sum_{i,j=1}^d ((w_n)_i - \mathbb{E}_T[(w_n)_i])x_i x_j ((w_n)_j - \mathbb{E}_T[(w_n)_j])\right] \\ &= \text{tr}(xx^T \text{Cov } w_n) = \langle x, (\text{Cov } w_n)x \rangle,\end{aligned}$$

## Theorem (Gauss Markov theorem)

*Suppose that the data obeys the linear model*

$$Y = \langle w, X \rangle + \epsilon,$$

*with  $\mathbb{E}[\epsilon|X = x] = 0$ ,  $\text{Var}[\epsilon|X = x] = \sigma^2$  and errors at different points are uncorrelated.*

*Then*

- *the least squares estimator  $w_n = (X^T X)^{-1} X^T Y$  is **unbiased**,*
- *among **all possible unbiased estimators** of the weight vector  $w$  it has the smallest variance.*



The Gauss-Markov-Theorem is only of very limited practical use:

- Model assumption has to be true !  
In reality linearity is not often encountered
- If the model assumption is correct:  
least squares estimator is the best among all possible **unbiased** estimators !  
 $\Rightarrow$  a slightly biased estimator (e.g. ridge regression or lasso) could have much smaller variance  
 $\Rightarrow$  better **expected squared error**  $\Rightarrow$  **Better estimator !**