

Linear Regression

Maschinelles Lernen 1 -
Grundverfahren WS20/21

Prof. Gerhard Neumann
KIT, Institut für Anthropomatik und Robotik

Learning Outcomes

- Get familiar with matrix computations and matrix calculus
- Understand the regression problem
- What do we mean with “linear representation”?
- Be able to derive the least squares solution
- What is the use of regularization in ridge regression?
- How to extend linear regression to non-linear function?

Today's Agenda!

Recap: Types of Machine Learning

Recap: Linear Algebra

- Vectors, Matrices and manipulation of those

Linear Regression:

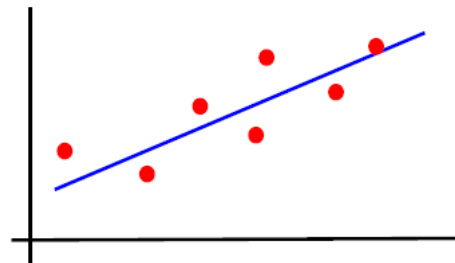
- Least-Squares Solution
- Generalized Linear Regression Models
- Ridge Regression

Supervised Learning

Training data includes targets

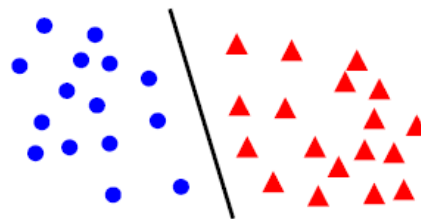
- **Regression:**

- Learn continuous function
- Example: line



- **Classification:**

- Learn class labels
- Example: Digit recognition

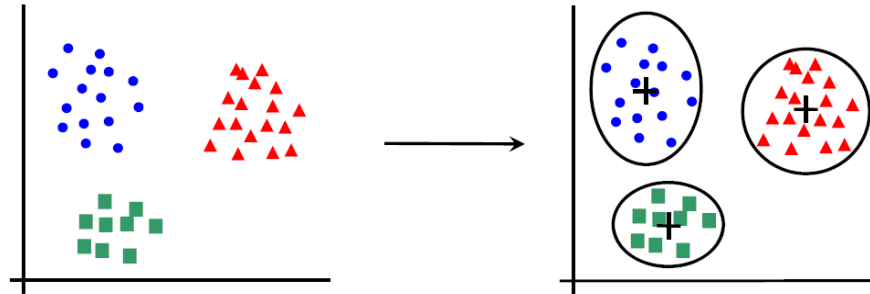


Unsupervised Learning

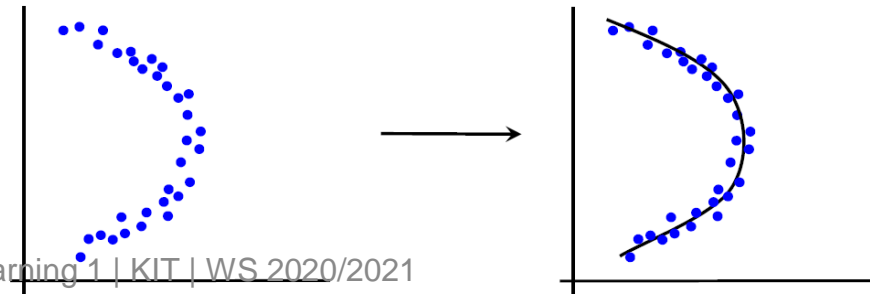
Trainings data does not include target values

- Model the data

- **Clustering:**



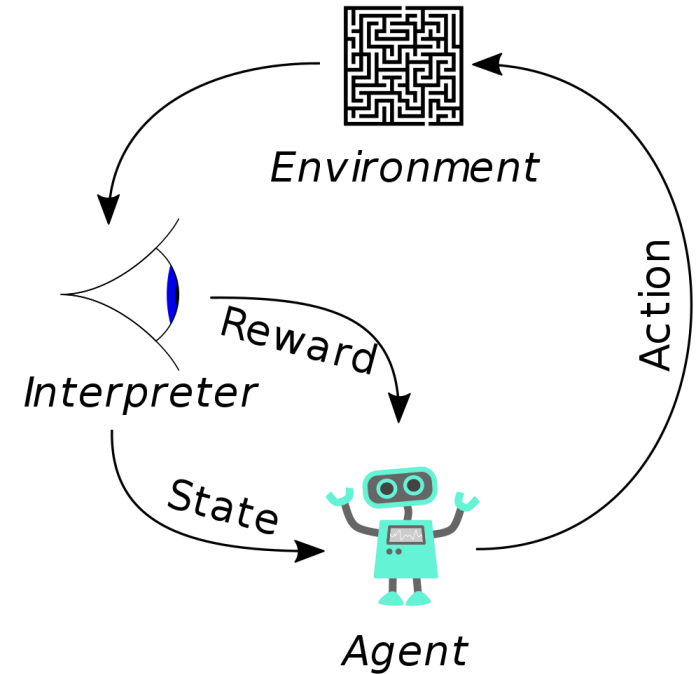
- **Dimensionality reduction:**



Reinforcement Learning

- No supervisor, but reward signal
- Selected actions also influence future states

Not part of this lecture!



Today's Agenda!

Recap: Types of Machine Learning

Recap: Linear Algebra

- Vectors, Matrices and manipulation of those

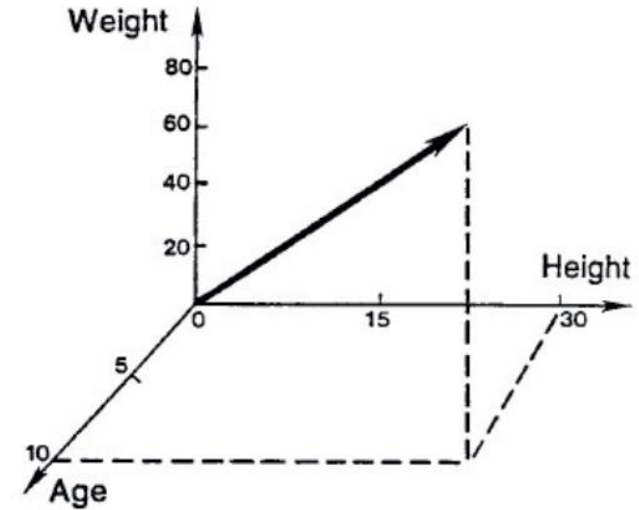
Linear Regression:

- Least-Squares Solution
- Generalized Linear Regression Models
- Ridge Regression

Vectors

- A vector is a multi-dimensional quantity

Joe	$\begin{bmatrix} 37 \\ 72 \\ 175 \end{bmatrix}$	Mary	$\begin{bmatrix} 10 \\ 30 \\ 61 \end{bmatrix}$
Carol	$\begin{bmatrix} 25 \\ 65 \\ 121 \end{bmatrix}$	Brad	$\begin{bmatrix} 66 \\ 67 \\ 155 \end{bmatrix}$



- Each dimension contains different information (Age, Height, Weight...)

Some notation

- Vectors will always be represented as bold symbols

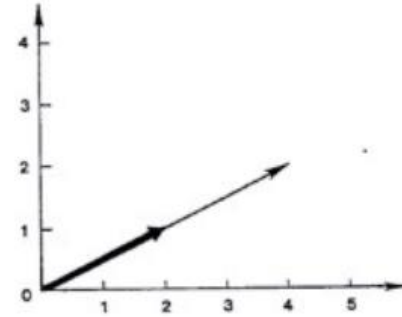
$$\underbrace{x = 1}_{\text{scalar}}, \quad \underbrace{\mathbf{x} = \begin{bmatrix} 1 \\ 2 \\ 4 \end{bmatrix}}_{\text{vector}}$$

- A vector \mathbf{x} is always a **column vector** $\mathbf{x} = \begin{bmatrix} 1 \\ 2 \\ 4 \end{bmatrix}$
- A transposed vector \mathbf{x}^T is always a **row vector** $\mathbf{x}^T = \begin{bmatrix} 1 & 2 & 4 \end{bmatrix}$

What can we do with vectors?

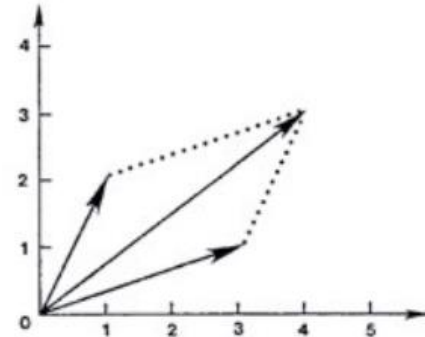
- Multiplication by scalars

$$2 \begin{bmatrix} 1 \\ 2 \\ 4 \end{bmatrix} = \begin{bmatrix} 2 \\ 4 \\ 8 \end{bmatrix}$$



- Addition of vectors

$$\begin{bmatrix} 1 \\ 2 \\ 4 \end{bmatrix} + \begin{bmatrix} 2 \\ 1 \\ 4 \end{bmatrix} = \begin{bmatrix} 3 \\ 3 \\ 8 \end{bmatrix}$$



Scalar products and length of vectors

- Scalar (Inner) products:
 - Sum the element-wise products

$$\boldsymbol{v} = \begin{bmatrix} 1 \\ 2 \\ 4 \end{bmatrix}, \quad \boldsymbol{w} = \begin{bmatrix} 2 \\ 4 \\ 8 \end{bmatrix}$$

$$\langle \boldsymbol{v}, \boldsymbol{w} \rangle = 1 \cdot 2 + 2 \cdot 4 + 4 \cdot 8 = 42$$

- Length of a vector
 - Square root of the inner product with itself

$$\|\boldsymbol{v}\| = \langle \boldsymbol{v}, \boldsymbol{v} \rangle^{\frac{1}{2}} = (1^2 + 2^2 + 4^2)^{\frac{1}{2}} = \sqrt{21}$$

Matrices

- A matrix is a rectangular array of numbers arranged in rows and columns.

$$\mathbf{X} = \begin{bmatrix} 1 & 3 \\ 2 & 3 \\ 4 & 7 \end{bmatrix} \quad \mathbf{A} = \begin{bmatrix} 1 & 3 & 5 & 4 \\ 2 & 3 & 7 & 2 \end{bmatrix}$$

- \mathbf{X} is a 3 x 2 matrix and \mathbf{A} a 2 x 4 matrix
- Dimension of a matrix is always *num rows times num columns*
- Matrices will be denoted with bold upper-case letters ($\mathbf{A}, \mathbf{B}, \mathbf{W}$)
- Vectors are **special cases of matrices**

$$\mathbf{x} = \underbrace{\begin{bmatrix} 1 \\ 2 \\ 4 \end{bmatrix}}_{3 \times 1 \text{ matrix}}$$

$$\mathbf{x}^T = \underbrace{\begin{bmatrix} 1 & 2 & 4 \end{bmatrix}}_{1 \times 3 \text{ matrix}}$$

Matrices in Machine Learning

- In many cases, our data set can be represented as matrix, where single samples are vectors

$$\text{Joe: } \mathbf{x}_1 = \begin{bmatrix} 37 \\ 72 \\ 175 \end{bmatrix} \quad \text{Mary: } \mathbf{x}_2 = \begin{bmatrix} 10 \\ 30 \\ 61 \end{bmatrix} \quad \text{Carol: } \mathbf{x}_3 = \begin{bmatrix} 25 \\ 65 \\ 121 \end{bmatrix} \quad \text{Brad: } \mathbf{x}_4 = \begin{bmatrix} 66 \\ 67 \\ 175 \end{bmatrix}$$

- Most typical representation:**

- Each row represent a data sample (e.g. Joe)
- Each column represents a data entry (e.g. age)

➡ \mathbf{X} is a *num samples x num entries* matrix

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \mathbf{x}_3^T \\ \mathbf{x}_4^T \end{bmatrix} = \begin{bmatrix} 37 & 72 & 175 \\ 10 & 30 & 61 \\ 25 & 65 & 121 \\ 66 & 67 & 175 \end{bmatrix}$$

What can you do with matrices?

- Multiplication with scalar

$$3\mathbf{M} = 3 \begin{bmatrix} 3 & 4 & 5 \\ 1 & 0 & 1 \end{bmatrix} = \begin{bmatrix} 9 & 12 & 15 \\ 3 & 0 & 3 \end{bmatrix}$$

- Addition of matrices

$$\mathbf{M} + \mathbf{N} = \begin{bmatrix} 3 & 4 & 5 \\ 1 & 0 & 1 \end{bmatrix} + \begin{bmatrix} 1 & 2 & 1 \\ 3 & 1 & 1 \end{bmatrix} = \begin{bmatrix} 4 & 6 & 6 \\ 4 & 1 & 2 \end{bmatrix}$$

- Matrices can also be transposed

$$\mathbf{M} = \begin{bmatrix} 3 & 4 & 5 \\ 1 & 0 & 1 \end{bmatrix}, \mathbf{M}^T = \begin{bmatrix} 3 & 1 \\ 4 & 0 \\ 5 & 1 \end{bmatrix}$$

Multiplication of a vector with a matrix

- Matrix-Vector Product: $\mathbf{u} = \mathbf{W}\mathbf{v} = \begin{bmatrix} 3 & 4 & 5 \\ 1 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 2 \end{bmatrix} = \begin{bmatrix} 3 \cdot 1 + 4 \cdot 0 + 5 \cdot 2 \\ 1 \cdot 1 + 0 \cdot 0 + 1 \cdot 2 \end{bmatrix} = \begin{bmatrix} 13 \\ 3 \end{bmatrix}$

- Think of it as:
$$\underbrace{\begin{bmatrix} \mathbf{w}_1, \dots, \mathbf{w}_n \end{bmatrix}}_{\mathbf{W}} \underbrace{\begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix}}_{\mathbf{v}} = \underbrace{\begin{bmatrix} v_1 \mathbf{w}_1 + \dots + v_n \mathbf{w}_n \end{bmatrix}}_{\mathbf{u}}$$

- Hence: $\mathbf{u} = v_1 \mathbf{w}_1 + \dots + v_n \mathbf{w}_n = 1 \begin{bmatrix} 3 \\ 1 \end{bmatrix} + 0 \begin{bmatrix} 4 \\ 0 \end{bmatrix} + 2 \begin{bmatrix} 5 \\ 1 \end{bmatrix} = \begin{bmatrix} 13 \\ 3 \end{bmatrix}$

- We sum over the columns \mathbf{w}_i of \mathbf{W} weighted by v_i

- Vector needs to have same dimensionality as number of columns!**

Multiplication of a matrix with a matrix

- Matrix-Matrix Product:

$$\mathbf{U} = \mathbf{W}\mathbf{V} = \begin{bmatrix} 3 & 4 & 5 \\ 1 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 3 \\ 2 & 4 \end{bmatrix} = \begin{bmatrix} 3 \cdot 1 + 4 \cdot 0 + 5 \cdot 2 & 3 \cdot 0 + 4 \cdot 3 + 5 \cdot 4 \\ 1 \cdot 1 + 0 \cdot 0 + 1 \cdot 2 & 1 \cdot 0 + 0 \cdot 3 + 1 \cdot 4 \end{bmatrix} = \begin{bmatrix} 13 & 32 \\ 3 & 4 \end{bmatrix}$$

- Think of it as: $\mathbf{W} \underbrace{\begin{bmatrix} \mathbf{v}_1, \dots, \mathbf{v}_n \end{bmatrix}}_{\mathbf{V}} = \begin{bmatrix} \underbrace{\mathbf{W}\mathbf{v}_1}_{\mathbf{u}_1}, \dots, \underbrace{\mathbf{W}\mathbf{v}_n}_{\mathbf{u}_n} \end{bmatrix} = \mathbf{U}$

- Hence: Each column $\mathbf{u}_i = \mathbf{W}\mathbf{v}_i$ in \mathbf{U} can be computed by a matrix-vector product

Multiplication of a matrix with a matrix

- **Dimensions:** $\underbrace{m \times n}_W \cdot \underbrace{n \times j}_V = \underbrace{m \times j}_U$
 - Number of columns of left matrix must match number of rows of right matrix
- **Non-commutative (in general):** $VW \neq WV$
- **Associative:** $V(WX) = (VW)X$
- **Transpose Product:** $(VW)^T = W^T V^T$

Important special cases

- **Scalar (Inner) product:**

$$\mathbf{w}^T \mathbf{v} = [w_1, \dots, w_n] \begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix} = w_1 v_1 + \dots + w_n v_n = \langle \mathbf{w}, \mathbf{v} \rangle$$

- The scalar product can be written as vector-vector product

Important special cases

- **Compute row/column averages of matrix** $\mathbf{X} = \underbrace{\begin{bmatrix} X_{1,1} & \dots & X_{1,m} \\ \vdots & & \vdots \\ X_{n,1} & \dots & X_{n,m} \end{bmatrix}}_{n \text{ (samples)} \times m \text{ (entries)}}$

- Vector of row averages (average over all entries per sample)

$$\begin{bmatrix} \frac{1}{m} \sum_{i=1}^m X_{1,i} \\ \vdots \\ \frac{1}{m} \sum_{i=1}^m X_{n,i} \end{bmatrix} = \mathbf{X} \begin{bmatrix} \frac{1}{m} \\ \vdots \\ \frac{1}{m} \end{bmatrix} = \mathbf{X} \mathbf{a}, \quad \text{with } \mathbf{a} = \begin{bmatrix} \frac{1}{m} \\ \vdots \\ \frac{1}{m} \end{bmatrix}$$

- Vector of column averages (average over all samples per entry)

$$\left[\frac{1}{n} \sum_{i=1}^n X_{i,1}, \dots, \frac{1}{n} \sum_{i=1}^n X_{i,m} \right] = \left[\frac{1}{n}, \dots, \frac{1}{n} \right] \mathbf{X} = \mathbf{b}^T \mathbf{X}, \quad \text{with } \mathbf{b} = \begin{bmatrix} \frac{1}{n} \\ \vdots \\ \frac{1}{n} \end{bmatrix}$$

Matrix Inverse

scalar

matrices

- **Definition:** $w \cdot w^{-1} = 1$ $WW^{-1} = I, \quad W^{-1}W = I$

- **Unit Element:** Identity matrix, e.g., 3 x 3: $I = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$

- **Verify it!**

$$W = \begin{bmatrix} 1 & \frac{1}{2} \\ -1 & 1 \end{bmatrix} \quad W^{-1} = \begin{bmatrix} \frac{2}{3} & -\frac{1}{3} \\ \frac{2}{3} & \frac{2}{3} \end{bmatrix}$$

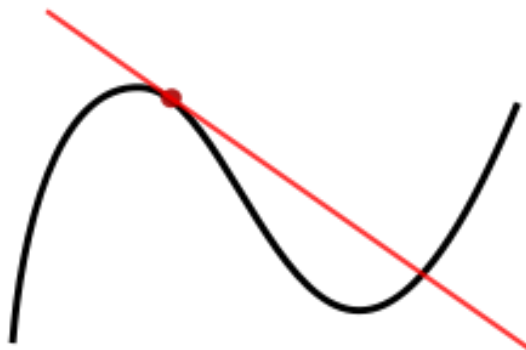
$$WW^{-1} = \begin{bmatrix} 1 & \frac{1}{2} \\ -1 & 1 \end{bmatrix} \begin{bmatrix} \frac{2}{3} & -\frac{1}{3} \\ \frac{2}{3} & \frac{2}{3} \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

- **Note:** We can only invert quadratic matrices (num rows = num cols)

Calculus

We also need to talk about derivatives...

“The derivative of a function of a real variable measures **the sensitivity to change of a quantity** (a function value or dependent variable) which is determined by another quantity (the independent variable)” (Wikipedia)

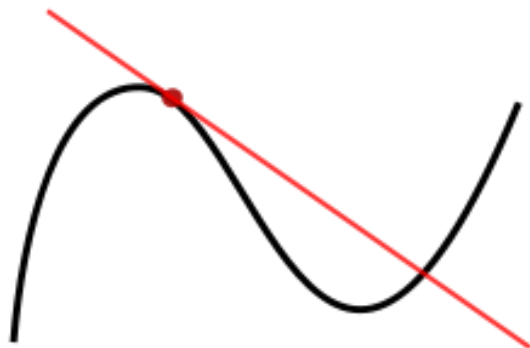


Function: $f(x)$

Derivative: $\frac{\partial f(x)}{\partial x}$

Minimum/Maximum: $\frac{\partial f(x)}{\partial x} = 0$

Derivatives and Gradients



Function:

scalar

$$f(x)$$

Derivative:

$$\frac{\partial f(x)}{\partial x} = g$$

Min/Max:

$$\frac{\partial f(x)}{\partial x} = 0$$

vector

$$f(\mathbf{x})$$

$$\frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} = \left[\frac{\partial f(\mathbf{x})}{\partial x_1}, \dots, \frac{\partial f(\mathbf{x})}{\partial x_d} \right]^T$$

$$\frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} = [0, \dots, 0]^T = \mathbf{0}$$

- $\frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} = \left[\frac{\partial f(\mathbf{x})}{\partial x_1}, \dots, \frac{\partial f(\mathbf{x})}{\partial x_d} \right]^T$ is called the gradient of function f at point \mathbf{x}
- We will use the „nabla“ operator as shorthand notation $\nabla f(\mathbf{x}) = \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}}$

Matrix Calculus

- We need to know some rules from **Matrix Calculus** (see wikipedia)

– Linear:

scalar

$$\frac{\partial ax}{\partial x} = a$$

– Quadratic:

$$\frac{\partial x^2}{\partial x} = 2x$$

vector

$$\nabla_x \mathbf{A} \mathbf{x} = \mathbf{A}^T$$

$$\nabla_x \mathbf{x}^T \mathbf{x} = 2\mathbf{x}$$

$$\nabla_x \mathbf{x}^T \mathbf{A} \mathbf{x} = 2\mathbf{A} \mathbf{x}$$

Today's Agenda!

Recap: Types of Machine Learning

Recap: Linear Algebra

- Vectors, Matrices and manipulation of those

Linear Regression:

- Least-Squares Solution
- Generalized Linear Regression Models
- Ridge Regression

Regression

Regression:

- Learn continuous function

$$y = f(x) + \epsilon$$

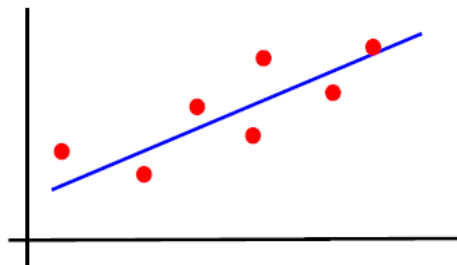
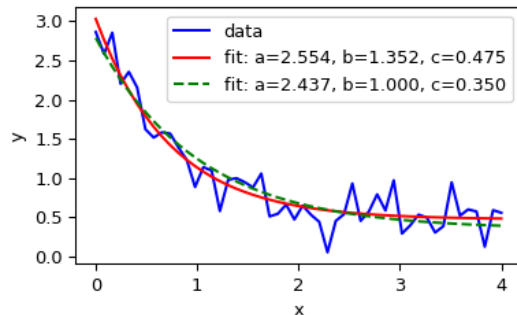
Linear Regression:

- We “just” fit a line

$$y = f(x) + \epsilon = w_0 + w_1x + \epsilon$$

We assume that the outputs are affected by (typically) Gaussian noise:

$$\epsilon \sim \mathcal{N}(0, 1)$$



Objective of Regression

We want to minimize the summed (or mean) squared error

$$\text{SSE} = \sum_{i=1}^N (y_i - f(\mathbf{x}_i))^2$$

- ... where the input \mathbf{x} is a d-dimensional vector

Why do we use the squared error?

- It is fully differentiable
- Easy to optimize
- It also makes sense as:

$$f^*(\mathbf{x}) = \operatorname{argmin}_{f(\mathbf{x})} \text{SSE} \Rightarrow f^*(\mathbf{x}) = \mathbb{E}[y|\mathbf{x}]$$

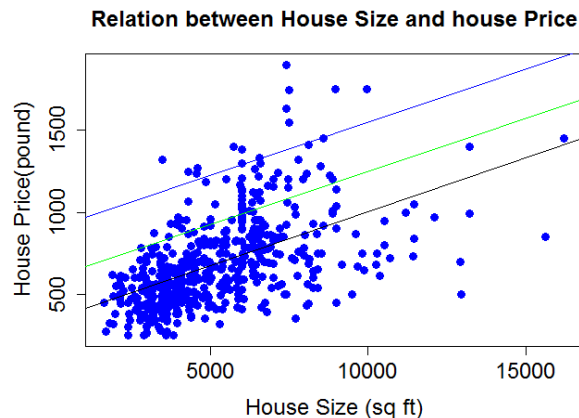
- Hence, we always estimate the mean of the target function!

Linear regression models

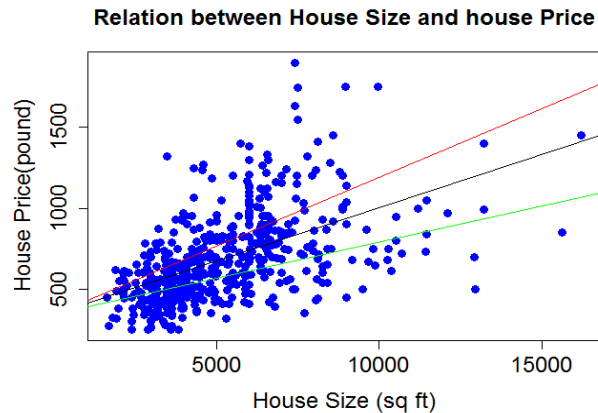
- In linear regression, the output y is modelled as linear function of the input \mathbf{x}_i

$$y = f(x) + \epsilon = w_0 + w_1x + \epsilon$$

Effect of w_0



Effect of w_1



Objective for Linear Regression

We want to consider linear functions with multiple inputs

$$f(\mathbf{x}_i) = w_0 + \sum_j w_j x_{i,j}$$

Our SSE objective now looks the following

$$\text{SSE} = \sum_{i=1}^N \left(y_i - \left(w_0 + \sum_j w_j x_{i,j} \right) \right)^2$$

Can we simplify it using matrices??

Linear regression models in matrix form

- Equation for the i-th sample

$$\hat{y}_i = w_0 + \sum_{j=1}^D w_j x_{i,j} = \tilde{\mathbf{x}}_i^T \mathbf{w}, \quad \text{with } \tilde{\mathbf{x}}_i = \begin{bmatrix} 1 \\ \mathbf{x}_i \end{bmatrix} \quad \text{and } \mathbf{w} = \begin{bmatrix} w_0 \\ \vdots \\ w_D \end{bmatrix}$$

- Equation for full data set

$$\hat{\mathbf{y}} = \begin{bmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_n \end{bmatrix} = \begin{bmatrix} \tilde{\mathbf{x}}_1^T \mathbf{w} \\ \vdots \\ \tilde{\mathbf{x}}_n^T \mathbf{w} \end{bmatrix} = \mathbf{X} \mathbf{w}$$

- $\hat{\mathbf{y}}$ is a vector containing the output for each sample

- $\mathbf{X} = \begin{bmatrix} \tilde{\mathbf{x}}_1^T \\ \vdots \\ \tilde{\mathbf{x}}_n^T \end{bmatrix} = \begin{bmatrix} 1 & \mathbf{x}_1^T \\ \vdots & \vdots \\ 1 & \mathbf{x}_n^T \end{bmatrix}$ is the data-matrix containing a vector of ones as the first column as bias

Linear regression models in matrix form

- Error vector:
$$\mathbf{e} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} - \begin{bmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_n \end{bmatrix} = \mathbf{y} - \hat{\mathbf{y}} = \mathbf{y} - \mathbf{X}\mathbf{w}$$

- Sum of squared errors (SSE)

$$\text{SSE} = \sum_i (y_i - \hat{y}_i)^2 = \sum_i e_i^2 = \mathbf{e}^T \mathbf{e} = (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w})$$

- We have now written the **SSE completely in matrix form!**

Deriving Linear Regression

- How do we obtain the optimal \mathbf{w} ? (which minimizes the SSE)

$$\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w}} \text{SSE} = \operatorname{argmin}_{\mathbf{w}} (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w})$$



At a minimal value of a function, its derivative is zero

I.e., find a \mathbf{w} where $\frac{\partial \text{SSE}}{\partial \mathbf{w}} = \mathbf{0}^T$

Estimation of \mathbf{w}

$$\begin{aligned}\text{SSE}(\mathbf{w}) &= (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w}) \\ &= \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{y}^T \mathbf{X} \mathbf{w} - \mathbf{w}^T \mathbf{X}^T \mathbf{y} + \mathbf{y}^T \mathbf{y} \\ &= \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} - 2\mathbf{y}^T \mathbf{X} \mathbf{w} + \mathbf{y}^T \mathbf{y}\end{aligned}$$

Take the derivative w.r.t \mathbf{w} :

$$\begin{aligned}\nabla_{\mathbf{w}} \text{SSE}(\mathbf{w}) &= \frac{\partial}{\partial \mathbf{w}} \left\{ \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} - 2\mathbf{y}^T \mathbf{X} \mathbf{w} + \mathbf{y}^T \mathbf{y} \right\} \\ &= \end{aligned}$$

Setting the gradient to $\mathbf{0}$ yields

$$\mathbf{w}^* =$$

Discussion

We have now derived our first ML algorithm: **Linear Regression!**

- The solution is called Least Squares solution
- One of the rare cases where we can obtain a closed form solution

This was only possible because:

- The cost-function (SSE) is convex for linear $f(x)$
 - There is only one minimum
- The cost function is quadratic in w
 - The minimum is easy to obtain

Ask questions!!!



Evaluating linear regression models

How can we estimate the quality of the model?

- The SSE can take arbitrary values depending on the range of the output
- Make the evaluation invariant to the variance of y

R-Square (or R^2) determines how much of the total variation in y is explained by the variation in x . Mathematically, it can be written as

$$R^2 = 1 - \frac{\text{Regression sum of squares}}{\text{Total sum of squares}} = 1 - \frac{\sum_{n=1}^N (\hat{y}_n - y_n)^2}{\sum_{n=1}^N (y_n - \bar{y})^2}$$

where \bar{y} is the mean of the outputs. R^2 tells how well the regression line approximates the real data points. An R^2 of 1 indicates that the regression line perfectly fits the data.

Today's Agenda!

Recap: Types of Machine Learning

Recap: Linear Algebra

- Vectors, Matrices and manipulation of those

Linear Regression:

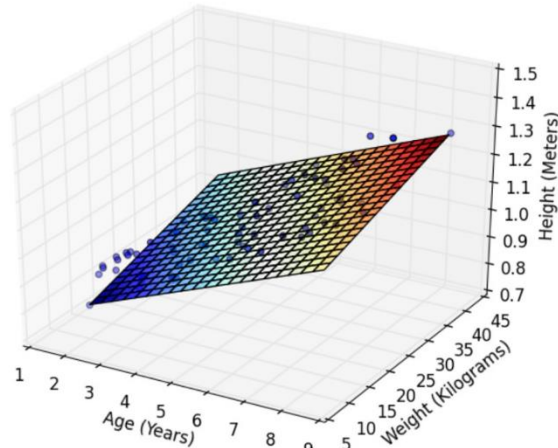
- Least-Squares Solution
- **Generalized Linear Regression Models**
- Ridge Regression

Linear Functions

So far, we modelled our function f as linear in \mathbf{x} and \mathbf{w}

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

However, this equation can only represent hyper-planes in the D -dimensional input space



General Form

In a more general writing, we could rewrite it as

$$f(x) = \phi(x)^T w$$

Where $\phi(x)$ is a vector valued function of the input vector x . This is also called **linear basis function models**, and $\phi_i(x)$ are known as **basis functions**.



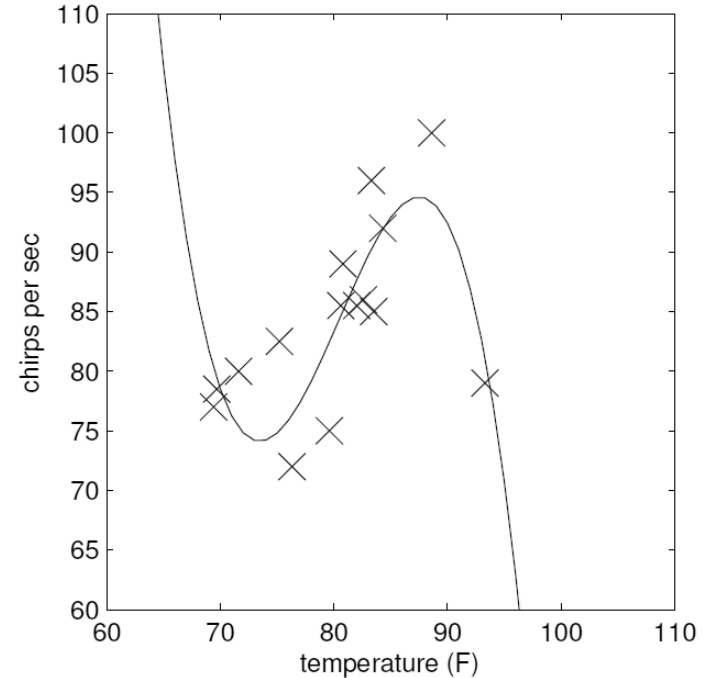
The model is linear in the parameter w , not necessarily linear in x .

Example of Polynomial Curve Fitting

$$f(\mathbf{x}) = \phi(\mathbf{x})^T \mathbf{w}$$

where

$$\mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ w_3 \end{bmatrix}, \quad \phi(\mathbf{x}) = \begin{bmatrix} 1 \\ x \\ x^2 \\ x^3 \end{bmatrix}$$

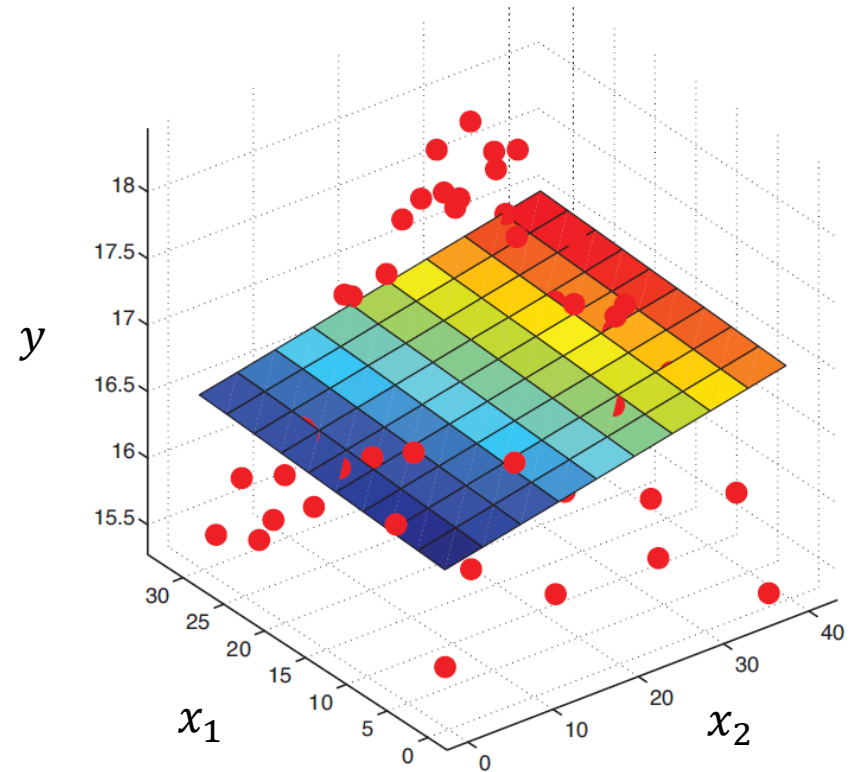


Example of Multiple Linear Regression

$$f(\mathbf{x}) = \phi(\mathbf{x})^T \mathbf{w}$$

where

$$\mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \end{bmatrix}, \quad \phi(\mathbf{x}) = \begin{bmatrix} 1 \\ x_1 \\ x_2 \end{bmatrix}$$

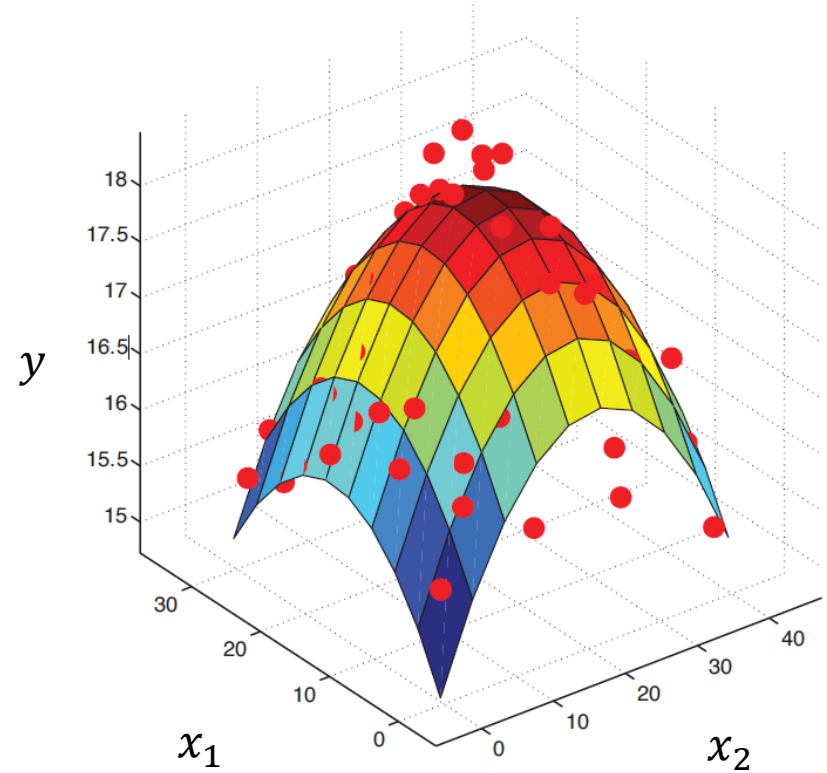


Example of Fitting Quadratic Form

$$f(\mathbf{x}) = \phi(\mathbf{x})^T \mathbf{w}$$

where

$$\mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ w_3 \\ w_4 \end{bmatrix}, \quad \phi(\mathbf{x}) = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ x_1^2 \\ x_2^2 \end{bmatrix}$$



Generalized Linear Regression

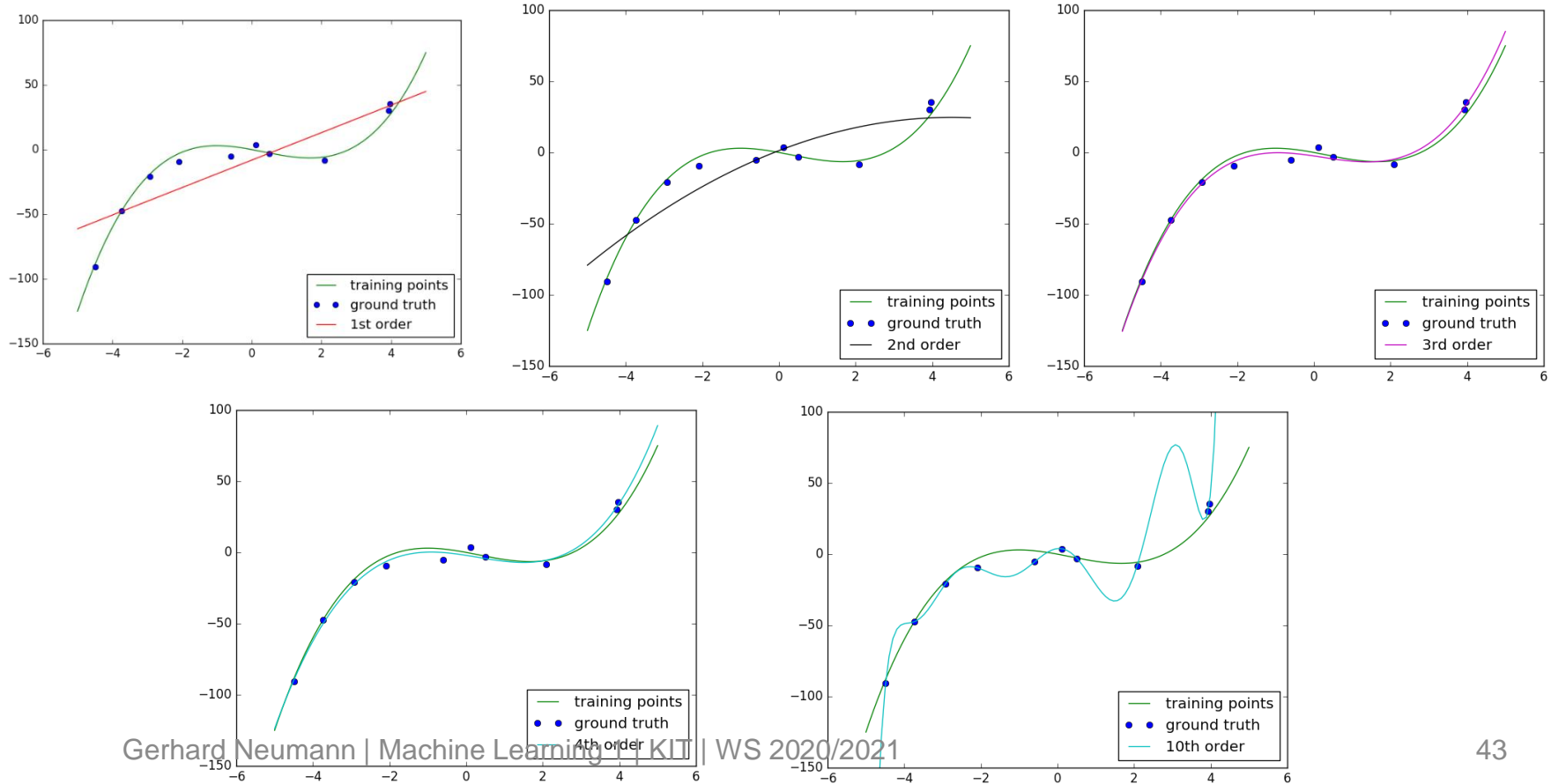
The derivations stay exactly the same, just the data matrix is now replaced by the basis function matrix, i.e.:

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

with $\Phi = \begin{bmatrix} \phi_1^T \\ \vdots \\ \phi_n^T \end{bmatrix}$

- In principle, this allows us to **learn any non-linear function**, if we know suitable basis functions (which is typically not the case).

Example: Selecting the order of the polynomial



Overfitting for polynomial regression

The error on the training set **is not an indication for a good fit!!**

- We always need an **independent test-set!**

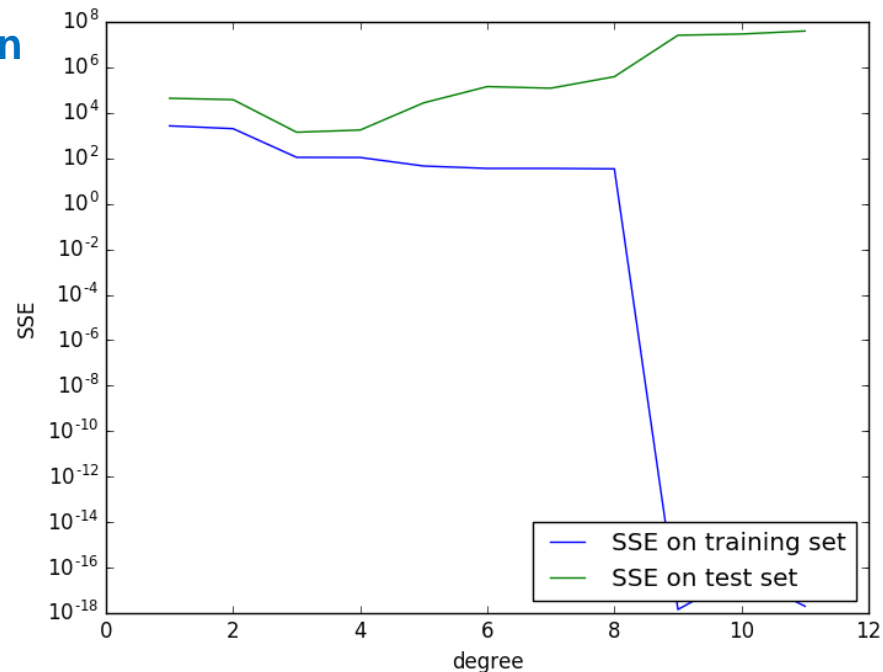
Overfitting:

- Training error goes down
- Test error goes up

The model is too complex. It fits the noise and has unspecified behavior between the training points.

Underfitting:

- Training + Test error are high
- The model is too simple to fit the data



Regularization

Regularization:

- Limit the model such that it can not fit the training data perfectly any more

Simple form of regularization: forcing the weights \mathbf{w} to be small

- Small weights will lead to a smoother function
- Introduce “regularization term” in cost-function

$$E_D(\mathbf{w}) + \lambda E_W(\mathbf{w})$$

Data term + Regularization term

- where λ is the regularization factor
- Needs to be tuned manually in most cases

Regularized Least Squares

With the sum-of-squares error function and a quadratic regularizer, we get

$$L_{\text{ridge}} = (\mathbf{y} - \Phi \mathbf{w})^T (\mathbf{y} - \Phi \mathbf{w}) + \lambda \mathbf{w}^T \mathbf{w}$$

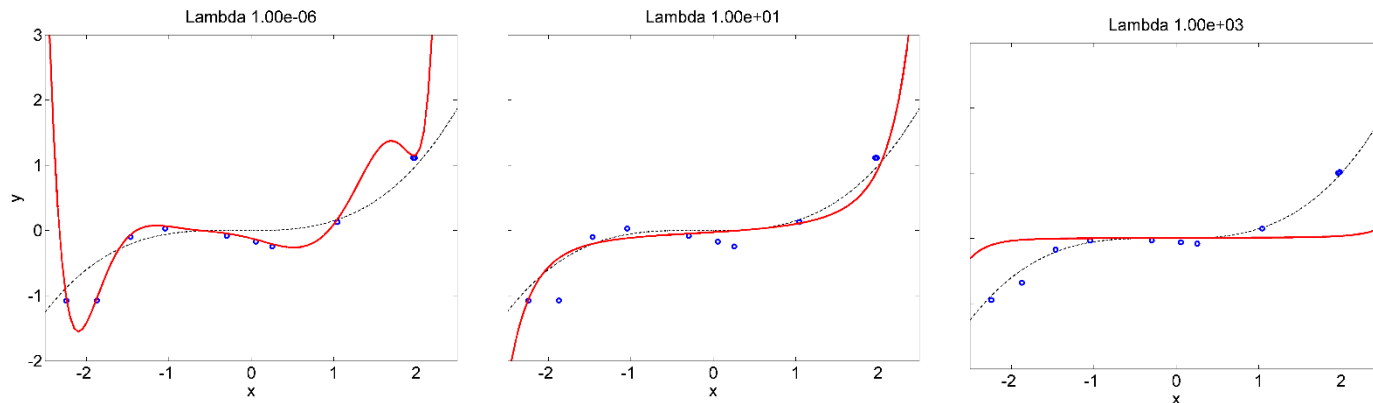
- This particular choice of regularizer is known as **weight decay** because in sequential learning algorithms, it encourages weight values to decay towards zero, unless supported by the data.
- In statistics, it is called **ridge regression**.

Derivations can be done similarly as before. The solution is given by

$$\mathbf{w}_{\text{ridge}}^* = (\Phi^T \Phi + \lambda \mathbf{I})^{-1} \Phi^T \mathbf{y}$$

- \mathbf{I} is the Identity matrix
- The matrix $(\Phi^T \Phi + \lambda \mathbf{I})$ is now full rank and can be more easily inverted

Ridge regression: Degree $n=15$



Influence of the regularization constant

Takeaway messages

What have we learned today?

- Familiarized with matrix manipulations and matrix calculus
- What a regression problem is
- How to obtain the Least-Squares solution in closed form
 - Only possible as the cost function is quadratic in the weights
- Generalized Linear Regression
 - Non-linear functions in x are fine as long as linear in w
- Avoid overfitting by keeping the weights small



