## **Model Linear untuk Regresi**

Dr. rer. nat. Hendri Murfi



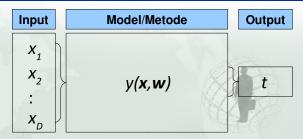
#### Intelligent Data Analysis (IDA) Group

Departemen Matematika, Universitas Indonesia – Depok 16424

Telp. +62-21-7862719/7863439, Fax. +62-21-7863439, Email. hendri@ui.ac.id

# Machine Learning

**Tahapan Umum Proses** 



Diberikan data pelatihan (training data), yaitu  $\mathbf{x}_i$  dan/atau  $\mathbf{t}_i$ , i = 1 sd N

- Preprocessing: pemilihan/ekstraksi fitur dari data, misal  $\mathbf{x}_i = (x_1, x_2, ..., x_p)^T$
- Learning: penentuan parameter metode, misal w, berdasarkan data pelatihan
- Testing: pengujian metode dengan data baru. Data penguji (testing data) tersebut harus dilakukan preprocessing yang sama dengan data pembelajaran sebelum dieksekusi oleh metode

### Learning

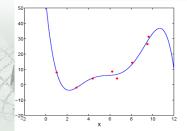
Diberikan data pelatihan  $\mathbf{x}_i$ , i = 1 sd N, dan/atau  $\mathbf{t}_i$ , i = 1 as N

- Supervised Learning. Data pelatihan disertai target, yaitu {x<sub>i</sub>, t<sub>i</sub>}, i = 1 sd
   N. Tujuan pembelajaran adalah membangun model yang dapat menghasilkan output yang benar untuk suatu data input, misal untuk regresi, klasifikasian, regresi ordinal, ranking, dll
- Unsupervised Learning. Data pelatihan tidak disertai target, yaitu x<sub>i</sub>, i = 1 sd N. Tujuan pembelajaran adalah membagun model yang dapat menemukan komponen/variabel/fitur tersembunyi pada data pelatihan, yang dapat digunakan untuk: pengelompokan (clustering), reduksi dimensi (dimension reduction), rekomendasi, dll

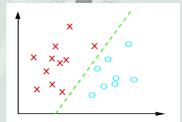
3

# Supervised Learning

- Regresi
  - Nilai output  $t_i$  bernilai kontinu (riil)
  - Bertujuan memprediksi output dengan akurat untuk data baru



- Klasifikasi
  - Nilai output  $t_i$  bernilai diskrit (kelas)
  - Bertujuan mengklasifikasi data baru dengan akurat



#### Regresi

#### **Model Linear**

 Model linear adalah kombinasi linear dari fungsi nonlinear dari variabel input (fungsi basis):

$$y(\mathbf{x}, \mathbf{w}) = \sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}) = \mathbf{w}^{\mathrm{T}} oldsymbol{\phi}(\mathbf{x})$$

dimana  $\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2, ..., \mathbf{x}_D)^{\mathsf{T}}$  adalah variabel input, dan  $\mathbf{w} = (\mathbf{w}_{0_i} \mathbf{w}_{1_i} ..., \mathbf{w}_{M-1})^{\mathsf{T}}$  adalah parameter,  $\phi(\mathbf{x}) = (\phi_0(\mathbf{x})_i \phi_1(\mathbf{x})_i ..., \phi_{M-1}(\mathbf{x}))^{\mathsf{T}}$  adalah vektor fungsi basis  $\phi_i(\mathbf{x})$ , M adalah jumlah total parameter dari model

- Biasanya,  $\phi_0(\mathbf{x}) = 1$ , sehingga  $\mathbf{w}_0$  berfungsi sebagai bias
- Ada banyak pilihan yang mungkin untuk fungsi basis  $\phi(\mathbf{x})$ , misal fungsi linear, fungsi polinomial, fungsi gaussian, fungsi sigmoidal, dll

5

# Regresi Linear Sederhana

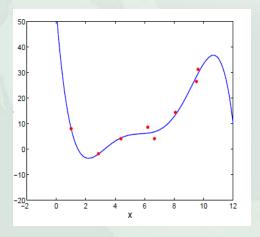
**Fungsi Basis Polinomial** 

Regresi linear sederhana (simple linear regression) adalah masalah regresi dengan variabel input x berdimensi satu. Misal kita menggunakan polinomial φ<sub>j</sub>(x) = x<sup>j</sup> sebagai fungsi basis, dan M = M-1, maka bentuk umum dari regresi linear sederhana tersebut adalah:

$$y(x, \mathbf{w}) = w_0 + w_1 x + w_2 x^2 + \ldots + w_M x^M = \sum_{j=0}^{M} w_j x^j$$

### Regresi Linear Sederhana

**Polynomial Curve Fitting** 



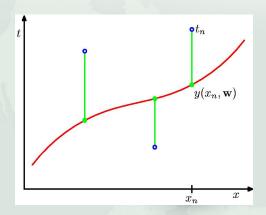
Diberikan data pelatihan  $\{x, t\}$ , i = 1 sd N

- Masalah: bagaimana mendapatkan kurva polinomial yang cocok untuk data pelatihan tersebut
- Solusi: mencari kurva polinomial yang memiliki kesalahan (error) terkecil pada data pelatihan tersebut
- Persoalan ini sering juga disebut sebagai polynomial curve fitting

7

# Regresi Linear Sederhana

Fungsi Error



 Salah satu fungsi error yang sering digunakan adalah fungsi sum-of-squares error sbb:

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^{N} \{y(x_n, \mathbf{w}) - t_n\}^2$$

 Salah satu metode yang digunakan untuk mencari nilai w yang meminimumkan fungsi error adalah metode kuadrat terkecil (least squares)

#### Regresi Linear Sederhana

Metode Kuadrat Terkecil

Setelah penurunan  $E(\mathbf{w})$  terhadap  $\mathbf{w}$ , maka persoalan penentuan nilai parameter w menjadi persoalan penentuan solusi sistem persamaan linear:

$$Aw = t$$

dimana

$$A = \begin{bmatrix} \sum_{n=1}^{N} 1 & \sum_{n=1}^{N} x_n & \cdots & \sum_{n=1}^{N} x_n^M \\ \sum_{n=1}^{N} x_n & \sum_{n=1}^{N} x_n^2 & \cdots & \sum_{n=1}^{N} x_n^{M+1} \\ \vdots & \vdots & \cdots & \vdots \\ \sum_{n=1}^{N} x_n^M & \sum_{n=1}^{N} x_n^{M+1} & \cdots & \sum_{n=1}^{N} x_n^{2M} \end{bmatrix} \qquad \mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_M \end{bmatrix} \qquad \mathbf{t} = \begin{bmatrix} \sum_{n=1}^{N} t_n \\ \sum_{n=1}^{N} x_n t_n \\ \vdots \\ \sum_{n=1}^{N} x_n^M t_n \end{bmatrix}$$

$$\mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_M \end{bmatrix}$$

$$\mathbf{t} = \begin{bmatrix} \sum_{n=1}^{N} t_n \\ \sum_{n=1}^{N} x_n t_n \\ \vdots \\ \sum_{n=1}^{N} x_n^M t_n \end{bmatrix}$$

# Regresi Linear Sederhana

**Contoh Kasus** 

Seorang ahli biologi telah melakukan eksperimen sebanyak 7 kali untuk melihat pertumbuhan bakteri berdasarkan kadar Nitrogen, dan diperoleh kondisi sbb:

Kadar Nitrogen (gram)	3	4	6	7	8	9
Pertumbuhan Bakteri	1	3	4	6	8	8

Tentukan regresi linear polinomial berorde 1 berdasarkan data tsb. Selanjutnya, prediksi pertumbuhan bakteri jika diberikan Nitrogen sebanyak 5 gram.

Solusi:

Dari persoalan diatas diketahui x = kadar nitrogen, t = pertumbuhan bakteri, N=6 dan M=1, sehingga:

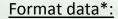
$$A = \begin{bmatrix} 6 & 37 \\ 37 & 255 \end{bmatrix}, \quad t = \begin{bmatrix} 30 \\ 217 \end{bmatrix}, \quad dan \ w = \begin{bmatrix} w_0 \\ w_1 \end{bmatrix} adalah \ solusi \ SPLAw = t \ , \ yaitu \ w = \begin{bmatrix} -2.35 \\ 1.19 \end{bmatrix}$$

dan model linear yang dihasilkan adalah y(x) = -2.35 + 1.19x. Sementara prediksi pertumbuhan bakteri untuk 5 gram Nitrogen adalah y(5) = -2.35 + 1.19\*5 = 3.6



# Regresi Linear Sederhana

Contoh Kasus: Menggunakan Weka



@RELATION bakteri

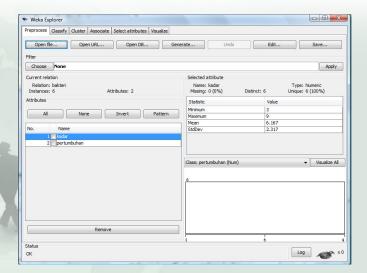
@ATTRIBUTE kadar NUMERIC
@ATTRIBUTE pertumbuhan NUMERIC

@DATA 3,1 4,3

7,6

8,8

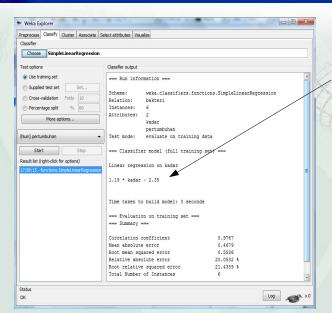
\*Disimpan dalam file dengan ekstensi arff (misal: bakteri.arff)



11

# Regresi Linear Sederhana

Contoh Kasus: Menggunakan Weka



Model hasil:

y(x) = -2.35 + 1.19x

#### Pemilihan Model

 Karakteristik model regresi linear polinomial ditentukan oleh nilai M (orde polinomial atau jumlah parameter).
 Pemilihan nilai M yang optimal dikenal juga dengan istilah pemilihan model (model selection)

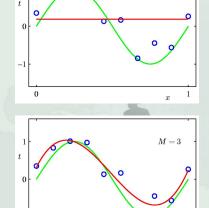
$$y(x, \mathbf{w}) = w_0 + w_1 x + w_2 x^2 + \ldots + w_M x^M = \sum_{j=0}^{M} w_j x^j$$

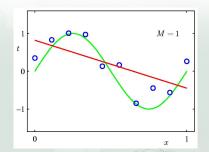
13

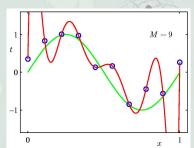
### Pemilihan Model

Under-fitting dan Over-fitting

M = 0

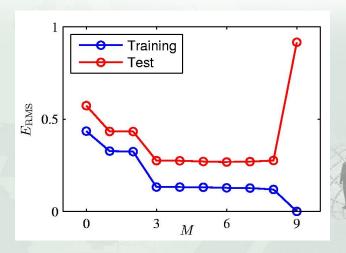








Under-fitting dan Over-fitting

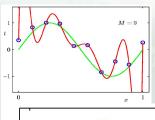


Root-Mean-Square (RMS) Error:

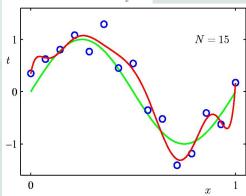
$$E_{
m RMS} = \sqrt{2E({f w}^{\star})/N}$$

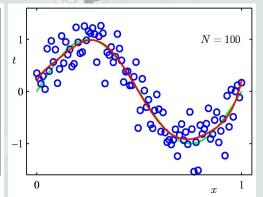
15

### Parameter vs Data



 Jumlah data pembelajaran seharusnya tidak lebih sedikit dari jumlah parameter





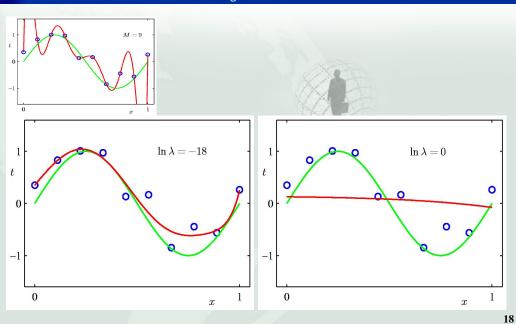
# Regularisasi

- Pada aplikasi praktis, kita sering menemukan kondisi dimana untuk persoalan yang kompleks ketersediaan data pembelajaran terbatas.
- Salah satu teknik yang digunakan untuk mengkontrol fenomena over-fitting adalah regularisasi (regularization), yaitu dengan cara menambah finalti ke fungsi error.

	M=0	M = 1	M=3	M = 9	, N
$w_0^\star$	0.19	0.82	0.31	0.35	$\simeq$ 1 $\sim$ 2 $\sim$ 32 $\wedge$ 10 $\sim$
$w_1^\star$		-1.27	7.99	232.37	$E(\mathbf{w}) = \frac{1}{2} \sum \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{2}{9} \ \mathbf{w}\ ^2$
$oldsymbol{w_2^{\star}}$			-25.43	-5321.83	2 n=1
$w_3^{ar{\star}}$			17.37	48568.31	
$w_4^\star$				-231639.30	
$w_5^{\star}$				640042.26	
$w_6^*$				-1061800.52	
$w_7^\star$				1042400.18	
$w_{8}^{\star}$				-557682.99	
$w_9^\star$				125201.43	17

Regularisasi

Penghalusan Kurva



# Regularisasi

Pengecilan Nilai Bobot

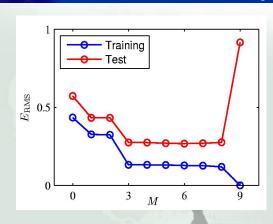
1	M=0	M = 1	M = 3	M = 9
$w_0^{\star}$	0.19	0.82	0.31	0.35
$w_1^\star$		-1.27	7.99	232.37
$w_2^\star$			-25.43	-5321.83
$w_3^\star$			17.37	48568.31
$w_4^\star$				-231639.30
$w_5^{\star}$				640042.26
$w_6^{\star}$	1			-1061800.52
$w_7^{\star}$	1/2		1	1042400.18
$w_8^\star$				-557682.99
$w_9^\star$	1			125201.43

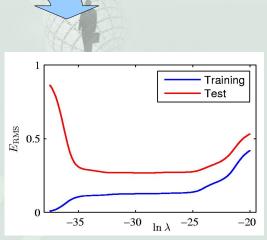
	$\ln \lambda = -\infty$	$\ln \lambda = -18$	$\ln \lambda = 0$
$w_0^{\star}$	0.35	0.35	0.13
$w_1^\star$	232.37	4.74	-0.05
$w_2^\star$	-5321.83	-0.77	-0.06
$w_3^\star$	48568.31	-31.97	-0.05
$w_4^\star$	-231639.30	-3.89	-0.03
$w_5^{\star}$	640042.26	55.28	-0.02
$w_6^{\star}$	-1061800.52	41.32	-0.01
$w_7^\star$	1042400.18	-45.95	-0.00
$w_8^\star$	-557682.99	-91.53	0.00
$w_9^\star$	125201.43	72.68	0.01

19

# Regularisasi

Mengatasi over-fitting





### Regresi Linear Umum

· Fungsi sum square error adalah

$$E_D(\mathbf{w}) = rac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^{\mathrm{T}} oldsymbol{\phi}(\mathbf{x}_n)\}^2$$

nilai bobot w yang meminimum fungsi error adalah

$$\mathbf{w}_{\mathrm{ML}} = \left(\mathbf{\Phi}^{\mathrm{T}}\mathbf{\Phi}
ight)^{-1}\mathbf{\Phi}^{\mathrm{T}}\mathbf{t}$$

dimana

$$oldsymbol{\Phi} = \left(egin{array}{cccc} \phi_0(\mathbf{x}_1) & \phi_1(\mathbf{x}_1) & \cdots & \phi_{M-1}(\mathbf{x}_1) \ \phi_0(\mathbf{x}_2) & \phi_1(\mathbf{x}_2) & \cdots & \phi_{M-1}(\mathbf{x}_2) \ dots & dots & \ddots & dots \ \phi_0(\mathbf{x}_N) & \phi_1(\mathbf{x}_N) & \cdots & \phi_{M-1}(\mathbf{x}_N) \end{array}
ight).$$

21

# Regresi Linear Umum

Regularisasi

· Fungsi regularized sum square error adalah

$$\frac{1}{2}\sum_{n=1}^{N}\{t_{n}-\mathbf{w}^{\mathrm{T}}\boldsymbol{\phi}(\mathbf{x}_{n})\}^{2}+\frac{\lambda}{2}\mathbf{w}^{\mathrm{T}}\mathbf{w}$$

nilai bobot w yang meminimum fungsi erroe adalah

$$\mathbf{w} = \left(\lambda \mathbf{I} + \mathbf{\Phi}^{\mathrm{T}} \mathbf{\Phi}\right)^{-1} \mathbf{\Phi}^{\mathrm{T}} \mathbf{t}.$$

dimana

$$oldsymbol{\Phi} = \left( egin{array}{cccc} \phi_0(\mathbf{x}_1) & \phi_1(\mathbf{x}_1) & \cdots & \phi_{M-1}(\mathbf{x}_1) \ \phi_0(\mathbf{x}_2) & \phi_1(\mathbf{x}_2) & \cdots & \phi_{M-1}(\mathbf{x}_2) \ dots & dots & \ddots & dots \ \phi_0(\mathbf{x}_N) & \phi_1(\mathbf{x}_N) & \cdots & \phi_{M-1}(\mathbf{x}_N) \end{array} 
ight).$$

