

# REGRESI LOGISTIK

Ali Akbar Septiandri

Universitas Al-Azhar Indonesia

*aliakbars@live.com*

April 5, 2020

# SELAYANG PANDANG

① ULASAN

② REGRESI LOGISTIK

③ OPTIMASI

④ KLASIFIKASI

# BAHAN BACAAN

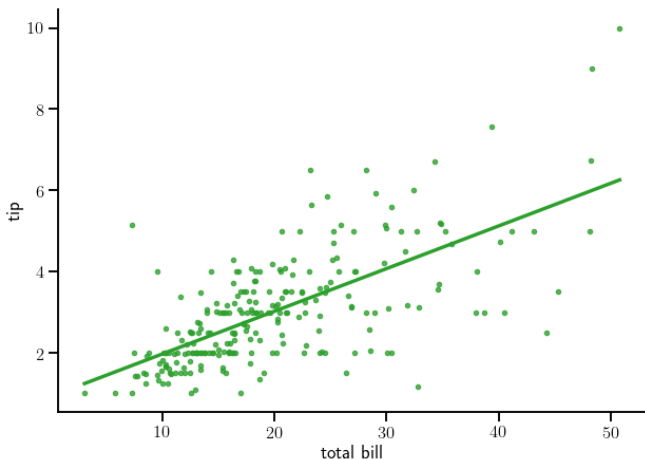
- 1 Septiandri, A.A. (2019). Artificial Intelligence Kuliah 3: Regresi Logistik. Github.
- 2 Murray, I. (2016). MLPR class notes. ([Regression and Gradients; Logistic Regression](http://www.inf.ed.ac.uk/teaching/courses/mlpr/2016/notes/)) <http://www.inf.ed.ac.uk/teaching/courses/mlpr/2016/notes/> (graduate level)

# ULASAN

## MINGGU LALU...

- Generalisasi error
- Bias-variance trade off
- Optimasi model: pembagian dataset
- Metrik evaluasi

# REGRESI LINEAR



**GAMBAR:** Mencari hubungan  $\mathbf{y} = X\mathbf{w}$  dengan meminimalkan  $E(\mathbf{w}) = \sum_{i=1}^N (y_i - \mathbf{w}^T \mathbf{x}_i)^2$

# MEMPREDIKSI KATEGORI

- Apa yang harus dilakukan jika kita ingin **memprediksi kategori** alih-alih *nilai riil*?

# MEMPREDIKSI KATEGORI

- Apa yang harus dilakukan jika kita ingin **memprediksi kategori** alih-alih *nilai riil*?
- Contoh: Prediksi apakah komentar-komentar berikut termasuk *spam* atau *ham* (bukan spam) jika dilihat dari kemunculan kata-kata 'order' dan 'password'.



# MEMPREDIKSI KATEGORI

- Apa yang harus dilakukan jika kita ingin **memprediksi kategori** alih-alih *nilai riil*?
- Contoh: Prediksi apakah komentar-komentar berikut termasuk *spam* atau *ham* (bukan spam) jika dilihat dari kemunculan kata-kata ‘order’ dan ‘password’.
- Kita asumsikan  $\text{spam} = 1$  dan  $\text{ham} = 0$ . Bagaimana memaksa keluaran dari regresi linear  $y \in (-\infty, \infty)$  menjadi  $y \in \{0, 1\}$ ?

# REGRESI LOGISTIK

## MENGUBAH KELUARAN

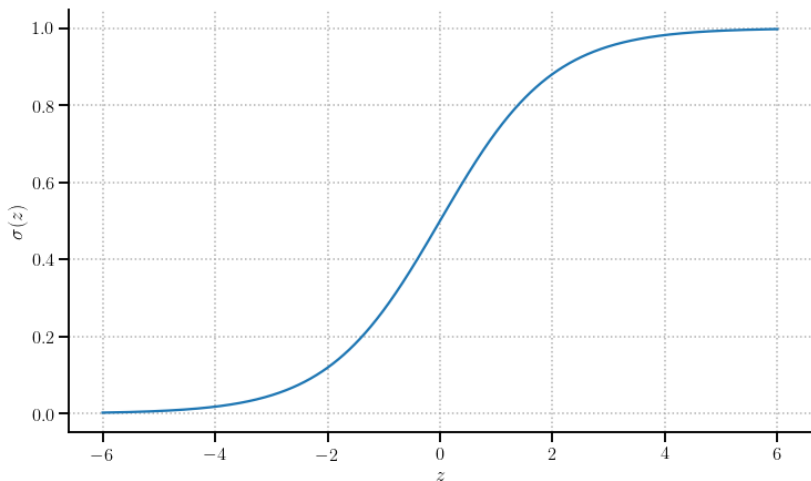
- Berdasarkan keluaran regresi linear, kita bisa memaksanya menjadi  $[0, 1]$
- Gunakan fungsi sigmoid:

$$P(y = 1|\mathbf{x}) = f(\mathbf{x}; \mathbf{w}) = \sigma(\mathbf{w}^T \mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}}$$

- Nilai  $[0, 1]$  dapat diartikan sebagai probabilitas dari kelas
- Karena probabilitas harus memiliki total 1, maka

$$P(y = 0|\mathbf{x}) = 1 - P(y = 1|\mathbf{x})$$

# FUNGSI SIGMOID



GAMBAR: Fungsi sigmoid/logistik  $\sigma(z) = \frac{1}{1+\exp(-z)}$

# DECISION BOUNDARY

- Dalam kasus satu variabel prediktor, kemiringan dari batas keputusan diatur oleh nilai  $w_1$ , sedangkan  $w_0$  (*intercept*) hanya menggesernya

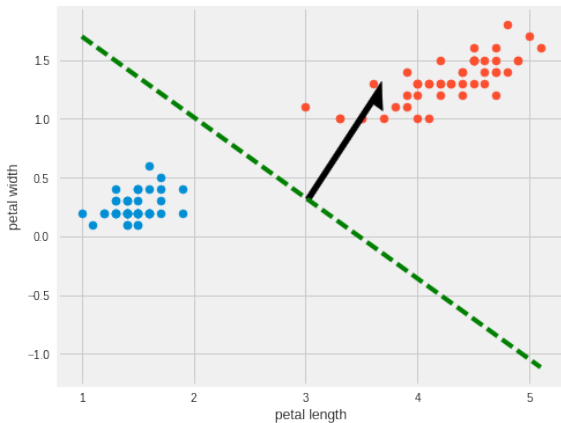
# DECISION BOUNDARY

- Dalam kasus satu variabel prediktor, kemiringan dari batas keputusan diatur oleh nilai  $w_1$ , sedangkan  $w_0$  (*intercept*) hanya menggesernya
- Batas keputusan yang dihasilkan akan berupa *hyperplane* yang akan tegak lurus terhadap vektor  $\mathbf{w}$

# DECISION BOUNDARY

- Dalam kasus satu variabel prediktor, kemiringan dari batas keputusan diatur oleh nilai  $w_1$ , sedangkan  $w_0$  (*intercept*) hanya menggesernya
- Batas keputusan yang dihasilkan akan berupa *hyperplane* yang akan tegak lurus terhadap vektor  $\mathbf{w}$
- Dari  $\mathbf{w}$ , kita bisa menggambarkan batas keputusan (*decision boundary*) ketika  $p(y = 1|\mathbf{x}) = p(y = 0|\mathbf{x}) = 0.5$ , i.e.  $\mathbf{w}^T \mathbf{x} = 0$

# DECISION BOUNDARY



**GAMBAR:** Batas keputusan dan vektor bobot untuk klasifikasi dua kelas



Bagaimana cara mencari nilai  $\mathbf{w}$ ?

## LIKELIHOOD

- Asumsi i.i.d.
- Dataset  $\mathcal{D} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\}$
- *Likelihood*-nya menjadi

$$\begin{aligned} p(\mathcal{D}|\mathbf{w}) &= \prod_{i=1}^N p(y = y_i|\mathbf{x}_i, \mathbf{w}) \\ &= \prod_{i=1}^N p(y = 1|\mathbf{x}_i, \mathbf{w})^{y_i} (1 - p(y = 1|\mathbf{x}_i, \mathbf{w}))^{1-y_i} \end{aligned}$$

- *Log likelihood*  $L(\mathbf{w}) = \log p(\mathcal{D}|\mathbf{w})$

$$L(\mathbf{w}) = \sum_{i=1}^N y_i \log \sigma(\mathbf{w}^T \mathbf{x}_i) + (1 - y_i) \log(1 - \sigma(\mathbf{w}^T \mathbf{x}_i))$$

## SOLUSI

- Nilai optimum untuk kasus ini unik, i.e. *convex*
- Untuk memaksimalkan nilainya, gunakan gradien

$$\frac{\partial L}{\partial w_j} = \sum_{i=1}^N (y_i - \sigma(\mathbf{w}^T \mathbf{x}_i)) x_{ij}$$

- Tidak ada solusi tertutup sehingga harus menggunakan *optimasi numerik*, e.g. dengan *gradient descent*

# OPTIMASI

Mengapa dinamakan *machine learning*?

# ALASAN MELAKUKAN OPTIMASI

- Belajar  $\rightarrow$  masalah optimasi kontinu
- Contoh: regresi linear, regresi logistik, jaringan saraf tiruan, SVM
- Salah satu caranya adalah dengan *maximum likelihood*

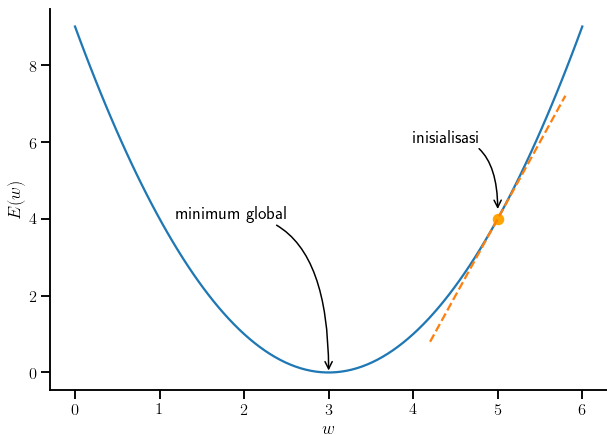
“Berapa peluangnya kita melihat data ini jika diketahui parameternya?”

# CARA MELAKUKAN OPTIMASI

- Menggunakan fungsi galat/error  $E(\mathbf{w})$  yang akan diminimalkan
- e.g. dapat berupa  $-L(\mathbf{w})$
- Beda nilai  $\mathbf{w}$ , beda besar error
- Belajar  $\equiv$  menuruni permukaan error



# MENURUNI PERMUKAAN FUNGSI ERROR



GAMBAR: Menuruni lembah fungsi error  $E(w)$

# GRADIENT DESCENT

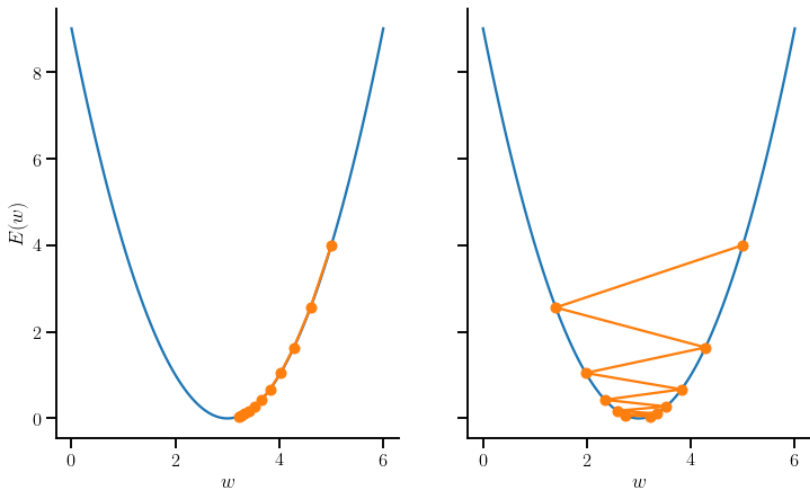
```
begin
|   Inisialisasi  $\mathbf{w}$ 
|   while  $E(\mathbf{w})$  masih terlalu besar do
|       |   Hitung  $\mathbf{g} \leftarrow \frac{\partial E}{\partial \mathbf{w}}$ 
|       |    $\mathbf{w} \leftarrow \mathbf{w} - \eta \mathbf{g}$ 
|   end
|   return  $\mathbf{w}$ 
end
```

**Algorithm 1:** Melatih dengan gradient descent

# LEARNING RATE

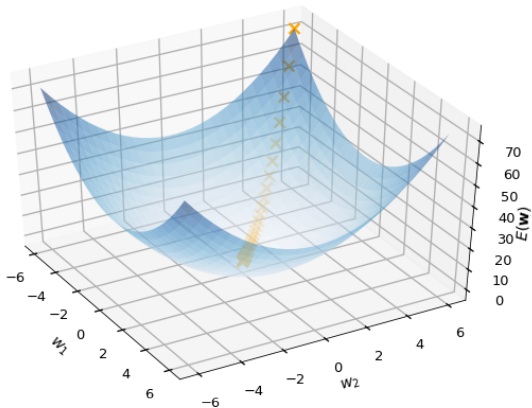
- $\eta$  (baca: “eta”) dikenal sebagai *step size* atau *learning rate* dengan nilai  $\eta > 0$
- $\eta$  terlalu kecil  $\rightarrow$  lambat
- $\eta$  terlalu besar  $\rightarrow$  tidak stabil

## PENGARUH NILAI $\eta$



**GAMBAR:** Fungsi  $E(w) = w^2 - 6w + 9$ , kiri:  $\eta = 0.1$ , kanan:  $\eta = 0.9$

# MENURUNI PERMUKAAN FUNGSI ERROR



GAMBAR: Menuruni lembah fungsi error

# BATCH VS ONLINE

- Untuk data yang sedikit, kita bisa menjumlahkan semua error sebelum memperbarui nilai  $\mathbf{w}$  (*batch*)

# BATCH VS ONLINE

- Untuk data yang sedikit, kita bisa menjumlahkan semua error sebelum memperbarui nilai  $\mathbf{w}$  (*batch*)
- Bagaimana untuk 10 juta data?

# BATCH VS ONLINE

- Untuk data yang sedikit, kita bisa menjumlahkan semua error sebelum memperbarui nilai  $\mathbf{w}$  (*batch*)
- Bagaimana untuk 10 juta data?
- Ternyata, kita bisa memperbarui nilai  $\mathbf{w}$  untuk setiap satu data (*online*)



# GRADIENT DESCENT (BATCH)

```
begin
  Inisialisasi  $\mathbf{w}$ 
  while  $E(\mathbf{w})$  masih terlalu besar do
    Hitung  $\mathbf{g} \leftarrow \sum_{i=1}^N \frac{\partial E_i}{\partial \mathbf{w}}$ 
     $\mathbf{w} \leftarrow \mathbf{w} - \eta \mathbf{g}$ 
  end
  return  $\mathbf{w}$ 
end
```

**Algorithm 2:** Melatih dengan batch gradient descent

# STOCHASTIC GRADIENT DESCENT

```
begin
  Inisialisasi  $\mathbf{w}$ 
  while  $E(\mathbf{w})$  masih terlalu besar do
    Pilih  $j$  sebagai integer acak antara 1..N
    Hitung  $\mathbf{g} \leftarrow \frac{\partial E_j}{\partial \mathbf{w}}$ 
     $\mathbf{w} \leftarrow \mathbf{w} - \eta \mathbf{g}$ 
  end
  return  $\mathbf{w}$ 
end
```

**Algorithm 3:** Stochastic gradient descent (SGD)

# KELEBIHAN DAN KEKURANGAN

- **Batch** lebih *powerful*
- **Batch** lebih mudah dianalisis
- **Online** lebih praktis untuk data yang besar

# PENGEMBANGAN GRADIENT DESCENT (NON-EXAMINABLE)

- “Why **Momentum** Really Works” [Goh, 2017]
- **Performance-dependent**  $\eta$ , e.g. “NewBOB”:  $\eta$  berubah menjadi setengahnya saat validation set tidak menjadi lebih baik
- **Time-dependent schedules**, e.g. eksponensial:  
 $\eta(t) = \eta(0)\exp(-t/r)$  ( $r \sim$  ukuran data latih)

## REGRESI LINEAR DENGAN GRADIENT DESCENT

[https://github.com/aliakbars/uai/blob/gh-pages/  
images/line.gif](https://github.com/aliakbars/uai/blob/gh-pages/images/line.gif)

# TENTANG METODE OPTIMASI

- Masih banyak metode optimasi yang tidak dibahas, e.g. linear programming, Newton's method, dll.
- Optimasi merupakan bidang matematika yang kompleks
- Masalah convex: optimum global. Non-convex: optimum lokal.
- Pahami mengapa *gradient descent* bisa mengalami masalah

# KLASIFIKASI

# MODEL GENERATIF DAN DISKRIMINATIF

- Naïve Bayes memodelkan bagaimana kelas “menghasilkan” vektor fitur  $p(\mathbf{x}|y)$  untuk kemudian diklasifikasikan dengan

$$p(y|\mathbf{x}) \propto p(\mathbf{x}|y)p(y)$$



# MODEL GENERATIF DAN DISKRIMINATIF

- Naïve Bayes memodelkan bagaimana kelas “menghasilkan” vektor fitur  $p(\mathbf{x}|y)$  untuk kemudian diklasifikasikan dengan

$$p(y|\mathbf{x}) \propto p(\mathbf{x}|y)p(y)$$

- Regresi logistik langsung memodelkan  $p(y|\mathbf{x})$ , i.e. diskriminatif

# MODEL GENERATIF DAN DISKRIMINATIF

- Naïve Bayes memodelkan bagaimana kelas “menghasilkan” vektor fitur  $p(\mathbf{x}|y)$  untuk kemudian diklasifikasikan dengan

$$p(y|\mathbf{x}) \propto p(\mathbf{x}|y)p(y)$$

- Regresi logistik langsung memodelkan  $p(y|\mathbf{x})$ , i.e. diskriminatif
- Keuntungan metode diskriminatif: Buat apa memodelkan  $p(\mathbf{x})$ ? Kita selalu punya input.

# MODEL GENERATIF DAN DISKRIMINATIF

- Naïve Bayes memodelkan bagaimana kelas “menghasilkan” vektor fitur  $p(\mathbf{x}|y)$  untuk kemudian diklasifikasikan dengan

$$p(y|\mathbf{x}) \propto p(\mathbf{x}|y)p(y)$$

- Regresi logistik langsung memodelkan  $p(y|\mathbf{x})$ , i.e. diskriminatif
- Keuntungan metode diskriminatif: Buat apa memodelkan  $p(\mathbf{x})$ ? Kita selalu punya input.
- Keuntungan metode generatif: Bisa menangani kasus data yang hilang, mendeteksi pencilan, atau mungkin *memang* perlu menghasilkan input

# KLASIFIKASI MULTIKELAS

- Buat vektor bobot  $\mathbf{w}_k$  untuk setiap kelas, untuk mengklasifikasikan  $k$  dan bukan- $k$
- Gunakan fungsi *softmax*

$$p(y = k|\mathbf{x}) = \frac{\exp(\mathbf{w}_k^T \mathbf{x})}{\sum_{j=1}^C \exp(\mathbf{w}_j^T \mathbf{x})}$$

- Perhatikan bahwa  $0 \leq p(y = k|\mathbf{x}) \leq 1$  dan  $\sum_{j=1}^C p(y = j|\mathbf{x}) = 1$

# IKHTISAR

- Klasifikasi dengan regresi logistik dan *gradient descent*
- Menggunakan  $-L(\mathbf{w})$  sebagai pengganti  $E(\mathbf{w})$
- Model generatif vs diskriminatif

# REFERENSI



Gabriel Goh (2017)

“Why Momentum Really Works”

Distill <http://distill.pub/2017/momentum/>

Terima kasih