

PERBANDINGAN ESTIMASI LASSO DAN LASSO *LEAST TRIMMED SQUARES* UNTUK MENGANALISIS DATA DIMENSI BESAR (Studi Kasus: Faktor – Faktor Yang Mempengaruhi Penyebaran Penyakit Tuberkulosis Di Sulawesi Selatan)

Comparison of LASSO and LASSO Least Trimmed Squares Estimation to Analyze High Dimensional Data (Case Study: Factors Affecting the Spread of Tuberculosis in South Sulawesi)

TRIGARCIA MALEACHI RANDA



**PROGRAM STUDI MAGISTER STATISTIKA
SEKOLAH PASCASARJANA
FAKULTAS MATEMATIKA DAN ILMU PENGETAHUAN ALAM
UNIVERSITAS HASANUDDIN
MAKASSAR
2022**

PERBANDINGAN ESTIMASI LASSO DAN LASSO *LEAST TRIMMED SQUARES* UNTUK MENGANALISIS DATA DIMENSI BESAR (Studi Kasus: Faktor – Faktor Yang Mempengaruhi Penyebaran Penyakit Tuberkulosis Di Sulawesi Selatan)

Tesis

sebagai salah satu syarat untuk mencapai gelar magister

Program Studi Magister Statistika

Disusun dan diajukan oleh

TRIGARCIA MALEACHI RANDA

H062202006

kepada

**PROGRAM STUDI MAGISTER STATISTIKA
FAKULTAS MATEMATIKA DAN ILMU PENGETAHUAN ALAM
UNIVERSITAS HASANUDDIN
MAKASSAR
2022**

TESIS

PERBANDINGAN ESTIMASI LASSO DAN LASSO *LEAST TRIMMED SQUARES* UNTUK MENGANALISIS DATA DIMENSI BESAR (Studi Kasus: Faktor – Faktor Yang Mempengaruhi Penyebaran Penyakit Tuberkulosis Di Sulawesi Selatan)

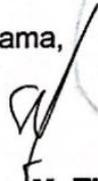
TRIGARCIA MALEACHI RANDA

H062202006

Telah dipertahankan di hadapan Panitia Ujian yang dibentuk dalam rangka Penyelesaian Program Studi Magister Statistika Fakultas Matematika dan Ilmu Pengetahuan Alam Universitas Hasanuddin pada tanggal 15 November 2022 dan dinyatakan telah memenuhi syarat kelulusan

Menyetujui,

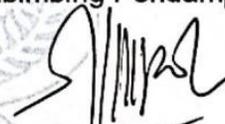
Pembimbing Utama,



Dr. Dr. Georgina M. Tinungki, M.Si.

NIP. 19620926 198702 2 001

Pembimbing Pendamping



Dr. Nurtiti Sunusi, S.Si M.Si.

NIP. 1972017 199703 2 002

Ketua Program Studi
Magister Statistika,



Dr. Dr. Georgina M. Tinungki, M.Si

NIP. 19620926 198702 2 001

Dekan Fakultas Matematika dan
Ilmu Pengetahuan Alam
Universitas Hasanuddin



Dr. Eng. Amiruddin, M.Si

NIP. 19720515 199702 1 002

PERNYATAAN KEASLIAN TESIS DAN PELIMPAHAN HAK CIPTA

Dengan ini saya menyatakan bahwa, tesis berjudul "Perbandingan Estimasi LASSO Dan LASSO Least Trimmed Squares Untuk Menganalisis Data Dimensi Besar (Studi Kasus: Faktor – Faktor Yang Mempengaruhi Penyebaran Penyakit Tuberkulosis Di Sulawesi Selatan)" adalah benar karya saya dengan arahan dari komisi pembimbing (Dr. Dr. Georgina Maria Tinungki, M.Si sebagai Pembimbing Utama dan Dr. Nurtiti Sunusi, S.Si., M.Si sebagai Pembimbing Pendamping). Karya ilmiah ini belum diajukan dan tidak sedang diajukan dalam bentuk apa pun kepada perguruan tinggi mana pun. Sumber informasinya yang berasal atau dikutip dari karya yang diterbitkan maupun tidak diterbitkan dari penulis lain telah disebutkan dalam teks dan dicantumkan dalam Daftar Pustaka tesis ini. Sebagian dari isi tesis ini telah dipublikasikan di Jurnal (EKSAKTA: Journal of Sciences and Data Analysis, ISSN: 2720-9326, Vol. 3, 103-112, DOI: 10.20885/EKSAKTA.vol3.iss2.art6) sebagai artikel dengan judul "Modeling the Proportion of Tuberculosis Cases in South Sulawesi using Sparse Least Trimmed Squares"

Dengan ini saya melimpahkan hak cipta dari karya tulis saya berupa tesis ini kepada Universitas Hasanuddin.

Makassar, 15 November 2022

Yang Menyatakan,



Handwritten signature of Trigarcia Maleachi Randa.

Trigarcia Maleachi Randa

NIM. H062202006

UCAPAN TERIMA KASIH

Puji syukur atas kehadiran Tuhan Yang Maha Esa atas limpahan rahmat dan karunia-Nya sehingga penulis dapat menyusun dan menyelesaikan tesis ini. Penulis menyadari sepenuhnya bahwa apa yang dikemukakan dalam tesis ini masih jauh dari kesempurnaan yang merupakan sebagai akibat dari keterbatasan kemampuan serta berbagai kesulitan yang penulis hadapi dalam menyusun tesis ini.

Penulis memanjatkan doa kepada Tuhan Yang Maha Esa agar memberikan rahmat-Nya kepada pihak yang banyak membantu dalam penyelesaian tesis ini. Penulis juga percaya tesis ini dapat selesai bukan hanya dengan kekuatan pikiran penulis semata akan tetapi karena bantuan dari berbagai pihak juga, baik selama proses perkuliahan bahkan sampai proses pengerjaan tesis di Program Magister Statistika Fakultas Matematika dan Ilmu Pengetahuan Alam, Universitas Hasanuddin. Namun demikian, penulis dengan senang hati menerima kritik dan saran yang bersifat membangun dari pembaca karya tulis ini demi sempurnanya tesis ini.

Terima kasih yang tak terhingga kepada kedua orang tua tercinta dan saudariku atas doa yang tak pernah putus, dukungan serta segala kebaikan mereka yang sampai kapan pun takkan pernah bisa terbalaskan atas kasih sayang yang tiada henti dalam penyelesaian tesis ini. Selanjutnya, saya ingin menyampaikan juga rasa hormat dan terima kasih kepada:

1. **Prof. Dr. Ir. Jamaluddin Jompa, M.Sc.** selaku Rektor Universitas Hasanuddin.
2. **Dr. Eng. Amiruddin, M.Si.** selaku Dekan Fakultas Matematika dan Ilmu Pengetahuan Alam beserta seluruh jajarannya.
3. **Dr. Nurtiti Sanusi, M.Si.** selaku Ketua Departemen Statistika yang menjadi Pembimbing Pertama yang telah bersabar dan bersedia meluangkan banyak waktunya untuk membimbing penulis dan memberikan masukan dalam penyelesaian tesis ini.
4. **Dr. Dr. Georgina Maria Tinungki, M.Si.** selaku Ketua Program Studi Magister Statistika yang menjadi Pembimbing Utama yang telah bersabar dan bersedia meluangkan banyak waktunya untuk membimbing penulis dan memberikan

ilmu, dukungan, dan motivasi serta kemudahan kepada penulis dalam berbagai hal selama menjalani Pendidikan di Departemen Statistika.

5. **Dr. Nirwan Ilyas, M.Si** selaku penguji penulis yang telah bersedia memberikan masukan-masukan dan arahan dalam penyusunan tesis.
6. **Dr. Anna Islamiyati, S.Si., M.Si** selaku penguji penulis yang telah bersedia memberikan masukan-masukan dan arahan dalam penyusunan tesis.
7. **Dr. Erna Tri Herdiani, M.Si** selaku penguji penulis yang telah bersedia memberikan masukan-masukan dan arahan dalam penyusunan tesis.
8. Teman - teman **Lab Statistika** yaitu **Agung Muhammad Takdir, Ratmila, Ainun Utari, Asnidar, Samsir Aditya Ania,** dan **Andi Dulung Laimbong** atas doa, semangat serta kebersamaannya selama ini yang banyak membantu penulis.
9. Teman – teman Mahasiswa Program Magister Statistika angkatan ketiga terkhusus terima kasih atas nasehat dan dukungan luar biasa kepada penulis, terkhusus **Kak Mubasiratul Munawaroh, A. Eka Hermia Fitriarningsy, Nurul Hidayanti Anggraini, Andi Isna Yunita, Alimatun Najiha, Maktisen Ena, Hedi Kuswanto, Muh. Qardawi Hamzah,** dan **Mushinil Haq.**

Semoga Allah SWT memberikan pahala yang berlipat ganda atas segala kebaikan yang telah diberikan kepada penulis dan semoga penulisan tesis ini bermanfaat bagi perkembangan ilmu pengetahuan dan teknologi, khususnya dalam dunia statistika dan data sains.

Penulis,

Trigarcia Maleachi Randa

ABSTRAK

TRIGARCIA MAELACHI RANDA. **Perbandingan Estimasi LASSO Dan LASSO Least Trimmed Squares Untuk Menganalisis Data Dimensi Besar (Studi Kasus: Faktor – Faktor Yang Mempengaruhi Penyebaran Penyakit Tuberkulosis Di Sulawesi Selatan)** (dibimbing oleh Georgina Maria Tinungki, dan Nurtiti Sunusi)

Tuberkulosis (TB) merupakan penyebab kematian ke-13 di dunia dan menjadi penyakit menular yang mematikan di Indonesia. Salah satu provinsi penyumbang kasus TB terbanyak di Indonesia pada tahun 2018 adalah Sulawesi Selatan dengan 84 kasus per 100,000 penduduk. Penelitian ini bertujuan untuk mengidentifikasi peubah yang dapat menjelaskan proporsi kasus TB di Sulawesi Selatan. Data yang digunakan memiliki banyak peubah bebas, dan terdapat pencilan. Analisis LASSO *Least Trimmed Squares* (LTS) dapat digunakan untuk menangani data yang memiliki banyak peubah bebas dan pencilan. Analisis data dilakukan pada data simulasi dan data proporsi kasus TB Berdasarkan rata-rata koefisien determinasi (R^2) dan *Root Mean Square Error* (RMSE), LASSO LTS memiliki performa terbaik dibandingkan LASSO pada setiap skenario kondisi ukuran data dan persentase pencilan. Analisis proporsi kasus TB di Sulawesi Selatan menghasilkan kesimpulan bahwa model LASSO LTS berhasil menyeleksi dan menyusutkan variabel menjadi 11 variabel dan menunjukkan hasil yang baik berdasarkan nilai R^2 dan RMSE untuk evaluasi model. Faktor-faktor tersebut bisa menjadi fokus pemerintah jika ingin menurunkan proporsi kasus TB di Sulawesi Selatan.

Kata Kunci: LASSO, pencilan, regresi terpenalti, regresi kekar, tuberkulosis

ABSTRACT

TRIGARCIA MAELACHI RANDA. **Comparison of LASSO and LASSO Least Trimmed Squares Estimation to Analyze High Dimensional Data (Case Study: Factors Affecting the Spread of Tuberculosis in South Sulawesi)** (supervised by Georgina Maria Tinungki, and Nurtiti Sunusi)

Tuberculosis is the 13th leading cause of death in the world and is a deadly infectious disease in Indonesia. One of the provinces that contributed the most TB cases in Indonesia in 2018 was South Sulawesi with 84 cases per 100,000 population. This study aims to identify variables that can explain the proportion of TB cases in South Sulawesi. The data used has many independent variables, and there are outliers. LASSO Least Trimmed Squares (LTS) analysis can be used for handle data that has many independent variables and outliers. Data analysis was performed on simulation data and proportion of TB cases data Based on the average coefficient of determination (R^2) and Root Mean Square Error (RMSE), LASSO LTS has the best performance than LASSO in each scenario condition of data size and outlier percentage. Analysis of proportion of TB cases in South Sulawesi resulted the conclusion that LASSO LTS model successfully select and shrink the variables to 11 variables only and shows satisfying results based on the value of R^2 and RMSE for the model evaluation. The government can focus on these factors if they want to reduce the proportion of TB cases in South Sulawesi.

Kata Kunci: LASSO, outliers, penalized regression, robust regression, tuberculosis

DAFTAR ISI

	Halaman
HALAMAN JUDUL	i
PERNYATAAN PENGAJUAN.....	ii
HALAMAN PENGESAHAN	iii
PERNYATAAN KEASLIAN TESIS	iv
UCAPAN TERIMA KASIH	v
ABSTRAK	vii
ABSTRACT	viii
DAFTAR ISI.....	ix
DAFTAR TABEL.....	xi
DAFTAR GAMBAR.....	i
DAFTAR LAMPIRAN	i
BAB I. PENDAHULUAN	1
1.1 Latar Belakang.....	1
1.2 Rumusan Masalah	3
1.3 Tujuan Penelitian	3
1.4 Manfaat Penelitian	4
1.5 Batasan Masalah	4
BAB II. TINJAUAN PUSTAKA	5
2.1 Analisis Regresi	5
2.2 Metode <i>Ordinary Least Square</i> (OLS)	6
2.3 Metode <i>Least Trimmed Squares</i> (LTS).....	6
2.4 <i>Least Absolute Shrinkage and Selection Operator</i> (LASSO).....	8
2.5 Algoritma <i>Least Angle Regression</i> (LAR).....	9
2.6 Validasi Silang (CV)	11
2.7 Estimasi Parameter Regresi Dengan Metode LASSO.....	12

2.8	Metode LASSO LTS.....	18
2.8.1	<i>Breakdown Point</i>	19
2.8.2	Algoritma LASSO LTS.....	20
2.8.3	<i>Reweighted LASSO LTS Estimator</i>	22
2.8.4	Pemilihan Parameter Penalti.....	24
2.9	Tuberkulosis	25
2.9.1	<i>Case Notification Rate (CNR)</i>	26
2.10	Kerangka Konseptual	29
BAB III. METODOLOGI PENELITIAN		30
3.1	Sumber Data.....	30
3.2	Variabel Penelitian	30
3.3	Struktur Data.....	35
3.4	Langkah-Langkah Penelitian.....	36
BAB IV. HASIL DAN PEMBAHASAN		39
4.1	Hasil Simulasi	39
4.2	Kajian Data Aktual.....	40
4.2.1	Eksplorasi Data	40
4.2.2	Pengujian Pencilan.....	42
4.2.3	Analisis LASSO.....	43
4.2.4	Analisis LASSO LTS.....	51
4.3	Pemilihan Model Terbaik	53
BAB V. KESIMPULAN DAN SARAN.....		54
5.1	Kesimpulan.....	54
5.2	Saran.....	55
DAFTAR PUSTAKA		56
LAMPIRAN		58

DAFTAR TABEL

Nomor Urut	Halaman
Tabel 1. Variabel Penelitian	30
Tabel 2. Struktur Data Penelitian	35
Tabel 3. Rataan nilai R^2 dan RMSE pada metode pada berbagai jumlah peubah dan persentase pencilan	39
Tabel 4. Plot korelasi <i>Pearson</i> antar peubah bebas	41
Tabel 5. Tahapan peubah bebas yang masuk ke dalam model	46
Tabel 6. Koefisien regresi hasil LASSO	50
Tabel 7. Koefisien regresi hasil LASSO LTS.....	52
Tabel 8. Ukuran kebaikan LASSO dan LASSO LTS.....	53

DAFTAR GAMBAR

Nomor Urut	Halaman
Gambar 1. Angka Notifikasi Semua Kasus TB Per 100,000 Penduduk Indonesia Tahun 2011-2020.....	26
Gambar 2. Angka Notifikasi Semua Kasus TB Per 100,000 Penduduk Menurut Provinsi Tahun 2018	27
Gambar 3. Angka Notifikasi Semua Kasus TB Per 100.000 Penduduk Provinsi Sulawesi Selatan Tahun 2012-2020.....	28
Gambar 4. Diagram Kerangka Konseptual Metode Penelitian	29
Gambar 5. Diagram alir pemodelan faktor-faktor yang mempengaruhi penyebaran penyakit TB di Indonesia.....	38
Gambar 6. Peta sebaran proporsi kasus TB di Sulawesi Selatan Tahun 2018	40
Gambar 7. Boxplot peubah respon (proporsi TB).....	43
Gambar 8. Plot pergerakan penduga koefisien regresi LASSO.....	45
Gambar 9. Plot validasi silang dengan menggunakan mode <i>step</i>	49
Gambar 10. Plot validasi silang dengan menggunakan mode <i>fraction</i>	50
Gambar 11. Validasi silang parameter LASSO LTS (λ)	51
Gambar 12. Plot antara residual standar dengan taksiran model	52

DAFTAR LAMPIRAN

Nomor Urut	Halaman
Lampiran 1. Nilai <i>p-value</i> antar peubah bebas	58
Lampiran 2. Koefisien regresi menggunakan metode LASSO dan nilai $\frac{\sum \beta_j }{\max \sum \beta_j }$ untuk setiap tahapan	60
Lampiran 3. Data proporsi kasus TB dan faktor-faktor yang mempengaruhi penyebaran penyakit TB di Sulawesi Selatan Tahun 2018	66

BAB I

PENDAHULUAN

1.1 Latar Belakang

Tuberkulosis (TB) adalah penyakit menular yang disebabkan oleh kuman *Mycobacterium tuberculosis* dan merupakan salah satu dari 10 penyebab utama kematian di seluruh dunia (WHO, 2018). Menurut *World Health Organization* (WHO) pada tahun 2018, Indonesia berada pada peringkat ketiga dengan penderita TB tertinggi di Dunia setelah India dan China. Sebanyak 10 juta orang menderita penyakit tuberkulosis, dan 1.2 juta meninggal karena penyakit ini pada tahun 2018. Sedangkan Provinsi Sulawesi Selatan digunakan dalam penelitian ini karena berdasarkan data dari Kementerian Kesehatan Republik Indonesia (Kemenkes RI) pada tahun 2018, Sulawesi Selatan merupakan salah satu Provinsi dengan jumlah semua kasus TB yang diobati dan dilaporkan di antara 100,000 penduduk tertinggi di Indonesia.

Kemudian pada tahun 2020 menurut data Kemenkes RI jumlah kasus tuberkulosis di Sulawesi Selatan sudah mengalami penurunan. Namun, penurunan jumlah kasus TB masih tetap harus diwaspadai karena adanya disparitas (kesenjangan) penyebaran penyakit TB antar daerah di Sulawesi Selatan, misalnya kabupaten Barru dan kabupaten Sidrap pada tahun 2018 memiliki jumlah kasus tuberkulosis berturut-turut sebanyak 182 dan 493 kasus, namun pada tahun 2020 mengalami kenaikan jumlah kasus secara berturut-turut sebesar 202 dan 591 kasus (BPS, 2020). Hal tersebut diduga karena datanya mengandung pencilan (*outlier*). Oleh sebab itu, perlu dilakukan penelitian lebih lanjut mengenai faktor-faktor yang mempengaruhi penyebaran penyakit tuberkulosis di Provinsi Sulawesi Selatan.

Seluruh data yang digunakan dalam penelitian ini diambil dari publikasi Riskesdas dan BPS pada tahun 2018. Peubah bebas yang dipilih adalah peubah yang berkaitan dengan kesehatan, ekonomi, sumber daya manusia (SDM) dan lingkungan. Faktor-faktor ini dipilih untuk diteliti karena menurut penelitian yang dilakukan oleh Sejati dan Sofiana (2015), secara keseluruhan penelitian tentang tuberkulosis melibatkan faktor-faktor tersebut. Keseluruhan data yang terkumpul memiliki 24 amatan, yaitu kabupaten/kota di Sulawesi Selatan, dengan 25 peubah

bebas yang terkait dengan kesehatan, ekonomi, SDM dan lingkungan. Peubah respon yang digunakan adalah proporsi tuberkulosis di Sulawesi Selatan tahun 2018 untuk masing-masing kabupaten/kota. Jika melihat perbandingan ukuran antara amatan dan jumlah peubah bebas yang digunakan, maka dapat dikatakan bahwa data ini termasuk yang memiliki dimensi besar, karena jumlah peubah bebasnya lebih besar dari jumlah amatannya.

Istilah dimensi besar mengacu pada kasus di mana banyaknya parameter yang tidak diketahui untuk diestimasi, p , jauh lebih besar daripada jumlah pengamatan, n , yaitu $p \gg n$. Saat memodelkan hubungan antara satu kejadian dengan kejadian lainnya, teknik statistika yang lazim digunakan adalah analisis regresi. Analisis regresi akan menghadapi masalah jika jumlah peubah bebas sangat banyak, sehingga diperlukan metode alternatif untuk menyelesaikan persoalan ini. Salah satu cara untuk menyelesaikan masalah ini di dalam literatur statistika dikenal metode regresi terpenalti.

Regresi terpenalti merupakan metode yang pendugaan parameternya didasarkan pada meminimumkan jumlah kuadrat galat yang diberi penalti (Hastie *et al.*, 2015). Metode regresi terpenalti yang sering digunakan untuk data dengan dimensi besar adalah regresi LASSO (*Least Absolute Shrinkage and Selection Operator*). LASSO merupakan metode komputasi dengan menggunakan pemrograman kuadrat yang dapat mengatasi masalah multikolinearitas dengan menyusutkan koefisien regresi mendekati nol bahkan hingga tepat nol sehingga dapat melakukan seleksi model (Tibshirani, 1996). Tetapi LASSO tidak bersifat *robust* terhadap pencilan (Alfons *et al.*, 2011). Terdapat beberapa metode regresi terpenalti yang bersifat *robust* terhadap pencilan dalam literatur dan di antaranya yang sangat direkomendasikan adalah LASSO *least trimmed squares* (LTS).

Menurut Wilems dan Aelst (2005), metode LTS merupakan suatu metode pendugaan parameter regresi *robust* untuk meminimumkan jumlah kuadrat sisaan sebanyak h . Akan tetapi, terdapat kekurangan pada metode LTS apabila data memiliki dimensi besar, maka metode ini tidak dapat dilakukan. Salah satu pengembangan metode LTS adalah LASSO LTS yang dapat digunakan mengatasi masalah ini. LASSO LTS diperkenalkan oleh Alfons *et al.* (2013), dimana metode ini merupakan gabungan antara metode LASSO dan LTS yang dapat mengatasi pencilan untuk data dimensi besar.

Beberapa penelitian sebelumnya mengenai metode LASSO dan LTS pernah dilakukan. Rohmawati *et al.* (2018) membandingkan metode estimasi LTS dan S (*Scale*) dalam mengestimasi model orde dua pada metode permukaan respon, dan diperoleh hasil estimasi LTS lebih baik dibandingkan estimasi S berdasarkan nilai R^2 . Yanke *et al.* (2022) meneliti tentang penanganan masalah multikolinearitas dalam pemodelan regresi pertumbuhan ekonomi Indonesia berdasarkan teori pertumbuhan ekonomi endogen. Dengan menerapkan metode seleksi peubah (*backward*, *forward*, dan *stepwise*), regresi komponen utama (PCR), *partial least square* (PLS), dan metode regularisasi (*Ridge*, LASSO, dan *Elastic Net*) untuk menyelesaikan masalah multikolinearitas, hasil yang diperoleh bahwa metode LASSO merupakan metode terbaik untuk mengatasi masalah multikolinearitas berdasarkan nilai *mean square error* (MSE).

1.2 Rumusan Masalah

Berdasarkan latar belakang yang telah diuraikan diatas, maka permasalahan utama yang ingin dibahas dalam penelitian ini adalah sebagai berikut:

1. Bagaimana perbandingan dari estimasi LASSO dan LASSO LTS?
2. Apa saja faktor-faktor yang mempengaruhi kasus penyebaran penyakit Tuberkulosis di Sulawesi Selatan berdasarkan model terbaik?

1.3 Tujuan Penelitian

Dari permasalahan diatas, maka tujuan penelitian ini adalah sebagai berikut:

1. Menunjukkan perbandingan estimasi LASSO dan LASSO LTS.
2. Mengetahui faktor-faktor yang mempengaruhi kasus penyebaran penyakit Tuberkulosis di Sulawesi Selatan berdasarkan model terbaik.

1.4 Manfaat Penelitian

Adapun manfaat yang ingin dicapai dalam penelitian ini adalah sebagai berikut:

1. Memberikan sumbangsih keilmuan dalam menerapkan estimasi LASSO dan LASSO LTS.
2. Sebagai bahan rujukan kepada Pemerintah dalam hal ini Kementerian Kesehatan, agar perlu memperhatikan faktor-faktor yang berpengaruh secara signifikan terhadap penyebaran penyakit tuberkulosis.

1.5 Batasan Masalah

Batasan masalah dalam penelitian ini adalah sebagai berikut:

1. Algoritma yang digunakan pada metode LTS adalah algoritma LTS yaitu penggabungan FAST-LTS dan *C-Steps*. Karena dalam proses estimasinya menggunakan algoritma LTS hanya akan memangkas sebaran data berdasarkan jumlah pencilan yang teramati sehingga akan menghasilkan fungsi objektif yang mengecil dan konvergen.
2. Algoritma yang digunakan pada metode LASSO adalah algoritma LARS yang merupakan modifikasi LAR (*Least Angle Regression*) untuk LASSO. Karena LARS merupakan algoritma yang lebih efisien digunakan dan telah dimodifikasi sehingga mempermudah dalam komputasi LASSO

BAB II

TINJAUAN PUSTAKA

2.1 Analisis Regresi

Analisis regresi merupakan suatu teknik statistika untuk memeriksa dan memodelkan hubungan antar peubah (Montgomery *et al.*, 2012). Analisis regresi merupakan teknik statistika yang paling banyak digunakan. Dalam analisis regresi terdapat istilah peubah Y dan peubah X. Peubah Y adalah peubah tak bebas atau peubah respon, sedangkan X adalah peubah bebas atau peubah prediktor. Analisis regresi yang hanya melibatkan satu peubah bebas disebut analisis regresi linier sederhana, sedangkan analisis regresi linier berganda melibatkan p peubah bebas. Hal utama dari analisis regresi adalah menduga parameter yang tidak diketahui dari model regresi.

Regresi linier sederhana memiliki beberapa asumsi, yaitu nilai harapan (rata-rata) sisaan sama dengan nol, ragam sisaan homogen, sisaan saling bebas, sisaan menyebar normal dengan rata-rata nol dan ragam (σ^2), serta sisaan bebas terhadap peubah bebas. Pada regresi linier berganda terdapat asumsi tambahan bahwa tidak ada multikolinearitas pada peubah bebas.

Model linier artinya linier dalam parameter. Model regresi linier memiliki bentuk sebagai berikut: (Montgomery dan Runger 2003)

$$y_i = \beta_0 + \sum_{j=1}^p \beta_j x_{ij} + e_i \quad (2.1)$$

dimana:

i : 1,2, ..., n

j : 1,2, ..., p

y_i : Vektor peubah respon berukuran $n \times 1$

x_{ij} : Matriks peubah bebas berukuran $n \times (p + 1)$

β_0 : Intersep

β_j : Slope atau kemiringan

e_i : Vektor sisaan (*error*) berukuran $n \times 1$

2.2 Metode Ordinary Least Square (OLS)

Metode kuadrat terkecil atau disebut juga *Ordinary Least Square* (OLS) merupakan metode yang digunakan untuk menduga koefisien regresi linier (β_0, β_j) dengan cara meminimumkan jumlah kuadrat sisaan (JKS) (Hastie *et al.*, 2008). Persamaan yang diminimumkan adalah sebagai berikut:

$$\text{JKS} = \sum_{j=1}^n \left(y_i - \beta_0 + \sum_{j=1}^p \beta_j x_{ij} + e_i \right)^2 \quad (2.2)$$

Penduga $\hat{\beta}$ dengan OLS akan menghasilkan penduga yang tak bias serta solusi unik. Pendugaan koefisien regresi dengan OLS memiliki kuadrat tengah sisaan terkecil di antara semua penduga linier yang tak bias. Namun, pada kondisi tertentu (misalnya multikolinieritas dan peubah sangat banyak), metode kuadrat terkecil sering tidak memuaskan. Hal tersebut disebabkan karena adanya masalah keakuratan prediksi yang mengakibatkan penduga kuadrat terkecil memiliki bias rendah tetapi ragam besar. Selain itu, semakin banyak peubah bebas maka model semakin sulit diinterpretasikan (Tibshirani, 1996).

2.3 Metode *Least Trimmed Squares* (LTS)

Menurut Chen (2002), estimasi *Least Trimmed Squares* (LTS) merupakan suatu metode *High Breakdown Value* yang diperkenalkan oleh Rousseeuw pada tahun 1984. Estimator LTS meminimumkan fungsi objektif yaitu jumlah kuadrat sisaan sebanyak h . LTS atau kuadrat terpangkas terkecil adalah salah satu teknik statistik yang digunakan untuk mengestimasi parameter dari model regresi linier yang memberikan alternatif *robust* ke metode regresi klasik berdasarkan meminimalkan jumlah kuadrat galat. Menurut Doornik (2011), penduga untuk LTS adalah sebagai berikut,

$$\hat{\beta}_{LTS} = \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^h e_{(i)}^2 \quad (2.3)$$

dimana,

$e_{(i)}^2$: kuadrat galat dari h pengamatan yang telah diurutkan $e_{(1)}^2 \leq e_{(2)}^2 \leq \dots \leq e_{(n)}^2$

h : konstanta pemotongan

Menurut Wilems dan Aelst (2005), h optimal yang digunakan pada metode LTS adalah

$$h = \left\lfloor \frac{n + p + 1}{2} \right\rfloor \quad (2.4)$$

Nilai h merupakan jumlah pengamatan yang digunakan untuk menduga parameter model regresi. Dengan p merupakan banyaknya peubah bebas dan n merupakan banyaknya pengamatan. Sama halnya dengan penduga lain pada regresi *robust*, pada metode LTS diberikan pembobot (w_{ii}) pada data sehingga data *outlier* tidak mempengaruhi model parameter hasil estimasi. Pembobot w_{ii} dapat disajikan dalam bentuk matriks, yang dinotasikan sebagai \mathbf{W} ,

$$\mathbf{W} = \begin{pmatrix} w_{11} & 0 & \cdots & 0 \\ 0 & w_{22} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & w_{nn} \end{pmatrix} \quad (2.5)$$

Menurut Chatzinakos dan Zioutas (2014), data pengamatan untuk metode LTS yang diidentifikasi sebagai pencilan/*outlier* diberikan bobot nol dan lainnya diberi bobot satu, dengan pembobotnya adalah sebagai berikut,

$$w_{ii} = \begin{cases} 1, & |e_i| \\ \hat{\sigma}^* & \text{lainnya} \end{cases} \quad (2.6)$$

dimana,

$$\hat{\sigma}^* = d_{h,n} \sqrt{\frac{1}{h} \sum_{i=1}^h e_{(i)}^2} \quad (2.7)$$

$$d_{h,n} = \frac{1}{\sqrt{1 - \frac{2n}{hc_{h,n}} \phi\left(\frac{1}{c_{h,n}}\right)}} \quad (2.8)$$

$$c_{h,n} = \frac{1}{\Phi^{-1}\left(\frac{h+n}{2n}\right)} \quad (2.9)$$

$$W_{scale} = \sqrt{\frac{\sum w_i e_i^2}{\sum w_i - k}} \quad (2.10)$$

Dengan ϕ merupakan fungsi kepadatan peluang dan distribusi normal standar, Φ merupakan fungsi distribusi kumulatif dan k merupakan banyaknya parameter pada model.

Menurut Willems and Aelst (2005) untuk menduga parameter model dengan menggunakan metode LTS digunakan algoritma LTS yaitu penggabungan *FAST-LTS* dan *C-Steps*, langkah-langkah algoritma LTS adalah sebagai berikut:

1. Menduga parameter model regresi menggunakan OLS
2. Menentukan kuadrat galat $e_{(i)}^2$, kemudian mengurutkan nilai $e_{(i)}^2$ dari yang terkecil ke terbesar
3. Menentukan nilai h yaitu pada persamaan (2.4) dan menghitung $\sum_{i=1}^h e_{(i)}^2$
4. Menghitung estimasi parameter $\hat{\beta}_{new}$ melalui OLS dari h pengamatan (nilai h awal)
5. Menentukan kuadrat residual $e_{(i) new}^2$ yang didapatkan dari estimasi parameter $\hat{\beta}_{new}$, kemudian menghitung h_{new} observasi dengan mengurutkan nilai $e_{(i) new}^2$ dari yang terkecil ke terbesar
6. Menghitung $\sum_{i=1}^{h_{new}} e_{(i) new}^2$
7. Ulangi langkah 4 sampai 6, iterasi dihentikan hingga dihasilkan hasil yang konvergen yaitu jika nilai $Q_{LTS}^{m+1} \leq Q_{LTS}^m$, m adalah jumlah iterasi
8. Memberikan bobot pada data dengan menerapkan *Final Weighted Least Square*.

2.4 Least Absolute Shrinkage and Selection Operator (LASSO)

Tibshirani pertama kali memperkenalkan metode Least Absolute Shrinkage and Selection Operator (LASSO) pada tahun 1996. Metode LASSO merupakan metode yang dapat digunakan untuk menyusutkan koefisien regresi dari peubah bebas yang memiliki korelasi tinggi dengan galat menjadi tepat pada nol atau mendekati nol. Sehingga dapat dikatakan bahwa metode LASSO berfungsi sebagai metode untuk melakukan seleksi variabel sekaligus hal tersebut dapat mengatasi masalah multikolinearitas (Tibshirani 1996). Penduga koefisien LASSO diperoleh dengan menggunakan pemrograman kuadratik karena tidak dapat diperoleh dalam bentuk tertutup seperti pada OLS atau regresi gulud. Metode LASSO mulai dikenal

setelah ditemukannya algoritma Least Angle Regression (LAR) pada tahun 2004 oleh Efron (Hastie *et al.*, 2008).

Menurut Zhao dan Yu (2006), persamaan secara umum LASSO dinyatakan sebagai berikut:

$$\mathbf{y}^{**} = \mathbf{X}^{**} + \boldsymbol{\beta} + \mathbf{e}^{**} \quad (2.11)$$

dimana:

\mathbf{y}^{**} : Vektor peubah respon berukuran $n \times 1$

\mathbf{X}^{**} : Matriks peubah bebas berukuran $(n \times p)$

$\boldsymbol{\beta}$: Vektor dari koefisien LASSO berukuran $(k + 1) \times 1$

\mathbf{e}^{**} : Vektor sisaan (*error*) berukuran $n \times 1$

Menurut Tibshirani (1996), penduga koefisien LASSO dapat diperoleh dengan meminimumkan jumlah kuadrat sisaan (JKS) dengan persamaan sebagai berikut:

$$\hat{\boldsymbol{\beta}}_{\text{LASSO}} = \underset{\boldsymbol{\beta}}{\operatorname{argmin}} \left\{ \sum_{j=1}^n \left(y_i = \beta_0 + \sum_{j=1}^p \beta_j x_{ij} + e_i \right)^2 \right\} \quad (2.12)$$

dengan syarat $\sum_{j=1}^k |\beta_j| \leq \lambda$. Nilai λ merupakan parameter *tuning* yang mengontrol penyusutan koefisien LASSO dengan $t \geq 0$.

Menurut Tibshirani (1996), jika $\lambda < \lambda_0$ dengan $\lambda_0 = \sum_{j=1}^p |\hat{\beta}_j|$ maka akan menyebabkan koefisien menyusut mendekati nol atau tepat pada nol atau tepat pada nol, sehingga LASSO akan berperan sebagai seleksi variabel. Akan tetapi jika $\lambda > \lambda_0$ maka penduga koefisien LASSO memberikan hasil yang sama dengan penduga kuadrat terkecil. Koefisien regresi LASSO ditentukan berdasarkan parameter *tuning* yang sudah dibakukan $s = \frac{\lambda}{\sum_{j=1}^k |\hat{\beta}_j^0|}$ dengan $\lambda = \sum_{j=1}^p |\hat{\beta}_j|$, $\hat{\beta}_j^0$ adalah penduga kuadrat terkecil untuk model penuh, nilai λ optimal diperoleh melalui validasi silang (Randa *et al.*, 2022).

2.5 Algoritma Least Angle Regression (LAR)

Least Angle Regression (LAR) merupakan suatu metode regresi yang algoritmanya dapat dimodifikasi menjadi algoritma komputasi untuk metode LASSO. LAR menghasilkan efisiensi algoritma dalam menduga koefisien LASSO dengan

komputasi yang lebih cepat dibandingkan pemrograman kuadratik Hastie *et al.* (2008).

LAR mempunyai kesamaan pada regresi dengan metode *forward stepwise*. Regresi dengan metode *forward stepwise* membentuk model secara bertahap dengan menambahkan sebuah peubah bebas pada satu waktu. Pada setiap tahapannya akan diidentifikasi peubah bebas terbaik untuk diikutsertakan ke dalam model dan dihitung kembali nilai penduga kuadrat terkecilnya. Least Angle Regression menggunakan cara kerja yang sama tetapi hanya memasukkan sebanyak peubah bebas sebagaimana mestinya. Pada langkah pertama dilakukan identifikasi peubah bebas yang paling berkorelasi dengan peubah respon. Berikut merupakan tahapan algoritma LAR asli: (Hastie *et al.*, 2008):

1. Membakukan peubah bebas yang digunakan sehingga memiliki nilai tengah nol dan ragam satu. Tujuan dilakukan pembakuan ini agar dapat membandingkan dugaan koefisien regresi yang memiliki ragam yang berbeda dalam suatu model. Dimulai dengan sisaan $\mathbf{r} = \mathbf{y} - \bar{y}, \beta_1, \dots, \beta_p = 0$
2. Mencari peubah bebas x_j yang paling berkorelasi dengan \mathbf{r} .
3. Mengubah nilai β_j yang pada awalnya bernilai 0 bergerak menuju koefisien kuadrat terkecil (x_j, \mathbf{r}) , sampai x_k yang lain memiliki korelasi yang kurang lebih sama dengan sisaan x_j .
4. Mengubah nilai β_j dan β_k bergerak ke arah koefisien gabungan kuadrat terkecil dengan sisaan sekarang dengan (x_j, x_k) , sampai suatu peubah bebas yang lain misal x_l memiliki korelasi yang cukup dengan sisaan (x_j, x_k) .

Modifikasi algoritma LAR untuk mendapatkan solusi LASSO adalah dengan memodifikasi langkah ke-4 menjadi:

Jika koefisien bukan nol mencapai nilai nol, keluarkan peubah tersebut dari gugus peubah aktif dan hitung kembali arah kuadrat terkecil bersama.

5. Mengulang langkah nomor 4 sampai semua p peubah bebas dimasukkan. Setelah $\min(N - 1, p)$ langkah, solusi model penuh untuk kuadrat terkecil diperoleh.

LAR selalu mengambil p (banyaknya peubah) langkah untuk mendapatkan penduga kuadrat terkecil secara penuh, sedangkan modifikasi LAR untuk LASSO dapat memiliki lebih dari p langkah untuk mendapatkannya. Algoritma LASSO

dengan memodifikasi LAR merupakan cara yang efisien dalam komputasi solusi masalah LASSO, terutama ketika jumlah peubah bebas yang digunakan jauh lebih banyak daripada data amatannya (Hastie *et al.*, 2008).

2.6 Validasi Silang (CV)

Salah satu metode pemilihan model terbaik bisa dilakukan menggunakan nilai validasi silang atau *Cross Validation* (CV). Menurut Hastie *et al.* (2008) validasi silang merupakan metode paling sederhana dan banyak dipakai secara luas untuk menduga galat prediksi. Idealnya, ketika data yang dimiliki memadai, akan dapat ditentukan suatu anak gugus data validasi dan digunakan untuk mengukur ketepatan model yang dimiliki. Namun seringkali data yang dimiliki terlalu sedikit sehingga tidak memungkinkan untuk dilakukan validasi secara langsung. Solusinya adalah melakukan validasi silang yang menggunakan sebagian data yang tersedia untuk membangun model dan sebagian data yang lain digunakan sebagai data pengujian model.

Validasi silang membagi data menjadi dua bagian, yaitu data training dan data test. Data training digunakan untuk membangun nilai $\hat{\beta}$ sedangkan data test digunakan untuk menguji kebaikan prediksi dari $X\hat{\beta}$. Nilai validasi silang yang diperoleh merupakan penduga bagi sisaan prediksi (Izenman, 2008). James *et al.* (2013) mendefinisikan data latih sebagai data pemodelan dan data uji sebagai data validasi. Pengelompokan data pemodelan dan data validasi dilakukan secara acak, dan setiap pengamatan memiliki peluang yang sama menjadi data validasi.

Salah satu metode tipe validasi silang adalah *k-fold*. Metode ini memiliki kelebihan ketika jumlah data amatan yang digunakan sedikit. Dalam validasi silang k-fold, semua observasi dipartisi secara acak ke dalam k subcontoh. Setiap sub-contoh digunakan sebagai data test dan sisanya digunakan sebagai data training. Proses validasi silang diulang sampai k kali, dan setiap satu sub-contoh digunakan hanya sekali dalam data test (Izenman, 2008).

Nilai kuadrat tengah sisaan hasil validasi silang (CV MSE) dihitung dengan menggunakan persamaan sebagai berikut:

$$CV\ MSE = \frac{1}{k} \sum_{k=1}^K \sum_{(x_i, y_i) \in T} (y_i - \hat{y}_{-k}(x_i))^2 \quad (2.13)$$

dengan $y_i - \hat{y}_{-k}(x_i)$ adalah dugaan y untuk x_i pada saat *fold* ke- k tidak digunakan dalam menduga model, dan y_i adalah elemen peubah respon ke- i pada data uji T . Menurut Izenman (2008), validasi silang yang sebaiknya digunakan adalah validasi silang *5-fold* atau *10-fold* karena menghasilkan nilai sisaan validasi silang dengan bias tinggi tetapi ragam rendah.

2.7 Estimasi Parameter Regresi Dengan Metode LASSO

Regresi digunakan untuk menduga nilai suatu respon dari peubah bebas yang sudah diketahui atau diasumsikan ada hubungan dengan respon. Nilai dugaan parameter dari model regresi linier dapat ditaksir dengan metode kuadrat terkecil atau *ordinary least squares* (OLS). OLS pada prinsipnya adalah meminimumkan jumlah kuadrat *error* (residual). Model umum persamaan regresi dengan sampel n dan jumlah prediktor p dinyatakan sebagai berikut

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e} \quad (2.14)$$

dimana

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}, \mathbf{X} = \begin{bmatrix} 1 & x_{11} & x_{21} & \cdots & x_{p1} \\ 1 & x_{12} & x_{22} & \cdots & x_{p2} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{1n} & x_{2n} & \cdots & x_{pn} \end{bmatrix}, \boldsymbol{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{bmatrix}, \boldsymbol{\varepsilon} = \begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{bmatrix} \quad (2.15)$$

Selanjutnya akan ditunjukkan penduga untuk β dengan meminimumkan jumlah kuadrat galat

$$\sum_{i=1}^n e_i^2 \quad (2.16)$$

Diketahui bahwa

$$e_i^2 = (y_i - \beta_0 - \beta_1 x_{1i} - \cdots - \beta_p x_{pi})^2 \quad (2.17)$$

maka turunan parsial terhadap $\beta_0, \beta_1, \dots, \beta_p$ dan menyamakan dengan nol sehingga diperoleh nilai estimasi model regresi linier adalah sebagai berikut:

$$\begin{aligned}
\frac{\partial s}{\partial \beta_0} &= -2 \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_{1i} - \dots - \beta_p x_{pi}) &= 0 \\
\frac{\partial s}{\partial \beta_1} &= -2 \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_{1i} - \dots - \beta_p x_{pi}) x_{1i} &= 0 \\
\frac{\partial s}{\partial \beta_2} &= -2 \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_{1i} - \beta_2 x_{2i} - \dots - \beta_p x_{pi}) x_{2i} &= 0 \\
&\vdots & \vdots \\
\frac{\partial s}{\partial \beta_p} &= -2 \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_{1i} - \dots - \beta_p x_{pi}) x_{pi} &= 0
\end{aligned} \tag{2.18}$$

Dengan demikian diperoleh persamaan berikut:

$$\begin{aligned}
n\hat{\beta}_0 + \hat{\beta}_1 \sum_{i=1}^n x_{1i} + \hat{\beta}_2 \sum_{i=1}^n x_{2i} + \dots + \hat{\beta}_p \sum_{i=1}^n x_{pi} &= \sum_{i=1}^n y_i \\
\hat{\beta}_0 \sum_{i=1}^n x_{1i} + \hat{\beta}_1 \sum_{i=1}^n x_{1i}^2 + \hat{\beta}_2 \sum_{i=1}^n x_{1i}x_{2i} + \dots + \hat{\beta}_p \sum_{i=1}^n x_{1i}x_{pi} &= \sum_{i=1}^n x_{1i}y_i \\
\hat{\beta}_0 \sum_{i=1}^n x_{2i} + \hat{\beta}_1 \sum_{i=1}^n x_{1i}x_{2i} + \hat{\beta}_2 \sum_{i=1}^n x_{2i}^2 + \dots + \hat{\beta}_p \sum_{i=1}^n x_{2i}x_{pi} &= \sum_{i=1}^n x_{2i}y_i \\
&\vdots \\
\hat{\beta}_0 \sum_{i=1}^n x_{pi} + \hat{\beta}_1 \sum_{i=1}^n x_{1i}x_{pi} + \hat{\beta}_2 \sum_{i=1}^n x_{2i}x_{pi} + \dots + \hat{\beta}_p \sum_{i=1}^n x_{pi}^2 &= \sum_{i=1}^n x_{pi}y_i
\end{aligned} \tag{2.19}$$

Berdasarkan persamaan (2.39) dapat dibuat dalam persamaan berikut:

$$\mathbf{X}^T \mathbf{X} \boldsymbol{\beta} = \mathbf{X}^T \mathbf{y} \tag{2.20}$$

untuk menyelesaikan persamaan (2.40), maka dikalikan kedua sisinya dengan invers dari $(\mathbf{X}^T \mathbf{X})$. Maka diperoleh penduga dari $\boldsymbol{\beta}$ adalah sebagai berikut:

$$\begin{aligned}
(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} \boldsymbol{\beta} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \\
\mathbf{I} \boldsymbol{\beta} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \\
\boldsymbol{\beta} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}
\end{aligned}$$

Sehingga solusi estimasi parameter regresi dengan OLS adalah

$$\hat{\boldsymbol{\beta}}_{\text{OLS}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \tag{2.21}$$

Suatu matriks khusus di mana invers-nya dapat diperoleh dengan mentransposkan disebut matriks orthogonal. Pada kasus matriks ortonormal dimana $\mathbf{X}^T\mathbf{X} = \mathbf{I}$, solusi estimasi parameter regresi LASSO memiliki solusi eksplisit dan dapat dengan mudah dibandingkan dengan solusi OLS. Jika diasumsikan tidak ada intersep persamaan dapat ditulis sebagai berikut

$$\begin{aligned}\hat{\boldsymbol{\beta}}_{\text{OLS}} &= (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y} \\ &= (\mathbf{I})^{-1}\mathbf{X}^T\mathbf{y} \\ &= \mathbf{X}^T\mathbf{y}\end{aligned}$$

Maka solusi OLS untuk kasus matriks ortonormal adalah

$$\hat{\boldsymbol{\beta}}_{\text{OLS}} = \mathbf{X}^T\mathbf{y} \quad (2.22)$$

Melalui bentuk Lagrangian, persamaan (2.12) dapat ditulis dalam persamaan berikut

$$\hat{\boldsymbol{\beta}}_{\text{LASSO}} = \underset{\boldsymbol{\beta}}{\operatorname{argmin}}(\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \lambda\|\boldsymbol{\beta}\|_1) \quad (2.23)$$

Kemudian persamaan (2.49) dapat diuraikan sebagai berikut

$$\begin{aligned}\hat{\boldsymbol{\beta}}_{\text{LASSO}} &= \underset{\boldsymbol{\beta}}{\operatorname{argmin}}\left(\left(\sqrt{(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^2}\right)^2 + \lambda|\boldsymbol{\beta}|\right) \\ &= \underset{\boldsymbol{\beta}}{\operatorname{argmin}}((\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^2 + \lambda|\boldsymbol{\beta}|) \\ &= \underset{\boldsymbol{\beta}}{\operatorname{argmin}}(\mathbf{y}^T\mathbf{y} - \mathbf{y}^T\mathbf{X}\boldsymbol{\beta} - \boldsymbol{\beta}^T\mathbf{X}^T\mathbf{y} + \boldsymbol{\beta}^T\mathbf{X}^T\mathbf{X}\boldsymbol{\beta} + \lambda|\boldsymbol{\beta}|) \\ &= \underset{\boldsymbol{\beta}}{\operatorname{argmin}}(\mathbf{y}^T\mathbf{y} - 2\boldsymbol{\beta}^T\mathbf{X}^T\mathbf{y} + \boldsymbol{\beta}^T\mathbf{I}\boldsymbol{\beta} + \lambda|\boldsymbol{\beta}|) \\ \hat{\boldsymbol{\beta}}_{\text{LASSO}} &= \underset{\boldsymbol{\beta}}{\operatorname{argmin}}(\mathbf{y}^T\mathbf{y} - 2\boldsymbol{\beta}^T\hat{\boldsymbol{\beta}}_{\text{OLS}} + \boldsymbol{\beta}^T\boldsymbol{\beta} + \lambda|\boldsymbol{\beta}|) \quad (2.24)\end{aligned}$$

Karena $\mathbf{y}^T\mathbf{y}$ tidak memiliki parameter $\boldsymbol{\beta}$, sehingga dapat dibuang dan diperoleh persamaan sebagai berikut

$$\begin{aligned}\hat{\boldsymbol{\beta}}_{\text{LASSO}} &= \underset{\boldsymbol{\beta}}{\operatorname{argmin}}(-2\boldsymbol{\beta}^T\hat{\boldsymbol{\beta}}_{\text{OLS}} + \boldsymbol{\beta}^T\boldsymbol{\beta} + \lambda|\boldsymbol{\beta}|) \\ \hat{\boldsymbol{\beta}}_{\text{LASSO}} &= \underset{\boldsymbol{\beta}}{\operatorname{argmin}}\sum_{j=1}^d(-2\boldsymbol{\beta}_j\hat{\boldsymbol{\beta}}_{\text{OLS}_j} + \boldsymbol{\beta}_j^2 + \lambda|\boldsymbol{\beta}_j|) \quad (2.25)\end{aligned}$$

Di mana X_j adalah kolom ke- j dari kolom matriks X . Dan karena terpisah maka dapat dipecahkan masalah optimasi secara terpisah untuk setiap peubah dalam penjumlahan. Kemudian untuk setiap $|\beta_j| = \beta$, sehingga perlu diselesaikan persamaan berikut

$$\hat{\beta}_{\text{LASSO}} = \underset{\beta}{\operatorname{argmin}}(-2\beta\hat{\beta}_{\text{OLS}} + \beta^2 + \lambda\beta) \quad (2.26)$$

Meskipun persamaan tersebut tidak dapat diturunkan, masalah ini dapat dipecahkan menjadi tiga kasus. Dalam kasus pertama, untuk $\hat{\beta}_{\text{OLS}} > 0$, maka solusi $\hat{\beta}_{\text{LASSO}} > 0$ sehingga perlu dicari solusi persamaan berikut

$$\underset{\beta > 0}{\operatorname{argmin}}(-2\beta\hat{\beta}_{\text{OLS}} + \beta^2 + \lambda\beta) \quad (2.27)$$

Langkah berikut dicari turunan terhadap β dan menyamakan dengan nol sehingga diperoleh

$$\begin{aligned} \frac{\partial}{\partial \beta}(-2\beta\hat{\beta}_{\text{OLS}} + \beta^2 + \lambda\beta) &= 0 \\ -2\hat{\beta}_{\text{OLS}} + 2\beta + \lambda &= 0 \\ 2\beta &= 2\hat{\beta}_{\text{OLS}} - \lambda \\ \beta &= \hat{\beta}_{\text{OLS}} - \frac{\lambda}{2} \end{aligned} \quad (2.29)$$

Dalam kasus kedua, untuk $\hat{\beta}_{\text{OLS}} < 0$, maka solusi $\hat{\beta}_{\text{LASSO}} < 0$ sehingga perlu dicari solusi persamaan berikut

$$\underset{\beta < 0}{\operatorname{argmin}}(-2\beta\hat{\beta}_{\text{OLS}} + \beta^2 - \lambda\beta) \quad (2.30)$$

Kemudian dicari turunan terhadap β dan menyamakan dengan nol sehingga diperoleh

$$\begin{aligned} \frac{\partial}{\partial \beta}(-2\beta\hat{\beta}_{\text{OLS}} + \beta^2 - \lambda\beta) &= 0 \\ -2\hat{\beta}_{\text{OLS}} + 2\beta - \lambda &= 0 \\ 2\beta &= 2\hat{\beta}_{\text{OLS}} + \lambda \\ \beta &= \hat{\beta}_{\text{OLS}} + \frac{\lambda}{2} \end{aligned} \quad (2.31)$$

Dalam kasus ketiga untuk $\hat{\beta}_{\text{OLS}} = 0$, maka solusi $\hat{\beta}_{\text{LASSO}} = 0$.

Jadi, ketika kolom matriks X ortogonal, solusi untuk estimasi parameter regresi LASSO adalah:

$$\hat{\boldsymbol{\beta}}_{\text{LASSO}} = \begin{cases} \hat{\boldsymbol{\beta}}_{\text{OLS}} - \frac{\lambda}{2}, & \text{untuk } \hat{\boldsymbol{\beta}}_{\text{OLS}} > 0 \\ \hat{\boldsymbol{\beta}}_{\text{OLS}} + \frac{\lambda}{2}, & \text{untuk } \hat{\boldsymbol{\beta}}_{\text{OLS}} < 0 \\ 0, & \text{lainnya} \end{cases} \quad (2.32)$$

Secara umum regresi LASSO tidak memiliki solusi secara eksplisit dalam menentukan koefisien taksirannya karena pada fungsi kendala regresi LASSO berbentuk fungsi mutlak yang tidak dapat diturunkan karena turunan kirinya berbeda dengan turunan kanannya, dengan demikian dibutuhkan pemrograman komputasi untuk menyelesaikannya. Algoritma LARS merupakan algoritma yang sangat efektif dalam membantu menyelesaikan solusi regresi LASSO secara komputasi.

Menurut Efron *et al.* (2004), LARS melakukan estimasi $\hat{\boldsymbol{\mu}} = \mathbf{X}^* \hat{\boldsymbol{\beta}}$, dengan langkah-langkah yang berurutan, dan di setiap langkah akan menambah satu kovariat ke dalam model. LARS merupakan modifikasi LAR untuk LASSO. Nilai $\hat{\boldsymbol{\mu}}$ diperoleh dari teknik iterasi dengan nilai awal $\hat{\boldsymbol{\mu}}^{(0)} = 0$. Langkah-langkah estimasi koefisien LASSO dengan algoritma LARS sebagai berikut:

Mencari vektor yang sebanding dengan vektor korelasi antara peubah prediktor dan galat dari setiap peubah bebas

$$\hat{\mathbf{c}} = \mathbf{X}^T(\mathbf{y} - \hat{\boldsymbol{\mu}}) \quad (2.33)$$

Menentukan korelasi saat mutlak terbesar

$$\hat{C} = \max\{|\hat{c}_j|\} \quad (2.34)$$

maka diperoleh $s_j = \text{sign}\{\hat{c}_j\}$ untuk $j \in A$

Menentukan X_A , Himpunan A merupakan himpunan indeks aktif dari peubah bebas $\{1, 2, 3, \dots, m\}$. Himpunan indeks aktif A ditentukan berdasarkan nilai korelasi mutlak terbesar. Didefinisikan matriks:

$$X_A = (\dots s_j X_j^* \dots)_{j \in A} \quad (2.35)$$

dengan X^* merupakan matrix X yang dinormalisasi dan tanda s_j bernilai ± 1 , maka

$$G_A = X_A^T X_A \text{ dan } P_A = (\mathbf{1}^T G_A^{-1} \mathbf{1}_A)^{-1/2} \quad (2.36)$$

Menghitung nilai vektor *equiangular*, vektor *equiangular* adalah suatu vektor yang membagi sudut dari kolom-kolom menjadi sama besar dengan besar sudutnya kurang dari 90° . Nilai vektor *equiangular* dicari menggunakan rumus sebagai berikut:

$$\mathbf{u}_A = X_A \omega_A \text{ dengan } \omega_A = P_A G_A^{-1} \mathbf{1}_A \quad (2.37)$$

Menghitung vektor *inner product*:

$$\mathbf{a} \equiv \mathbf{X}^T \mathbf{u}_A \quad (2.38)$$

Menghitung $\hat{\boldsymbol{\mu}}_A$ dengan $\hat{\boldsymbol{\mu}}_{A^+} = \hat{\boldsymbol{\mu}}_A + \hat{\boldsymbol{\gamma}} \mathbf{u}_A$

$$\hat{\boldsymbol{\gamma}} = \min_{j \in A^c}^+ \left\{ \frac{\hat{C} - \hat{c}_j}{P_A - a_j}, \frac{\hat{C} - \hat{c}_j}{P_A - a_j} \right\} \quad (2.39)$$

$\min_{j \in A^c}^+$ menunjukkan bahwa yang dipilih adalah nilai minimum positif dari j yang bukan merupakan himpunan A . Pada tahap akhir dalam memperoleh nilai $\hat{\boldsymbol{\gamma}}$ menggunakan rumus

$$\hat{\boldsymbol{\gamma}}_m = \frac{\hat{C}_m}{A_m} \quad (2.40)$$

Tanda dari koordinat bukan nol $\hat{\beta}_j$ sama dengan tanda \hat{c}_j , dengan rumus (Efron *et al.*, 2004):

$$\text{sign}(\hat{\beta}_j) = \text{sign}(\hat{c}_j) = s_j \quad (2.41)$$

didefinisikan bahwa $\hat{d} = s_j \omega_{A_j}$ sehingga $\hat{\boldsymbol{\mu}}$ menjadi:

$$\boldsymbol{\mu}(\boldsymbol{\gamma}) = \mathbf{X}\boldsymbol{\beta}(\boldsymbol{\gamma}) \text{ dengan } \beta_j(\boldsymbol{\gamma}) = \hat{\beta}_j + \gamma \hat{d}_j \quad (2.42)$$

$\beta_j(\boldsymbol{\gamma})$ akan berubah tanda pada saat

$$\gamma_j = \frac{-\hat{\beta}_j}{\hat{d}_j} \quad (2.43)$$

perubahan pertama yang terjadi pada

$$\tilde{\boldsymbol{\gamma}} = \min_{\gamma_j > 0} \{\gamma_j\} \quad (2.44)$$

Jika $\tilde{\gamma} < \hat{\gamma}$ maka $\beta_j(\gamma)$ bukan merupakan solusi LASSO karena pada $\beta_j(\gamma)$ telah berubah tanda melainkan pada c_j tidak berubah tanda dan proses LARS berhenti pada $\gamma = \tilde{\gamma}$ dan menghapus j dari perhitungan vektor *equiangular* berikutnya dan peubah j dimasukkan kembali pada tahap perhitungan LARS selanjutnya, maka:

$$\hat{\boldsymbol{\mu}}_{A_+} = \hat{\boldsymbol{\mu}}_A + \hat{\gamma} \mathbf{u}_A \text{ dan } A_+ = A - \{j\} \quad (2.45)$$

2.8 Metode LASSO LTS

Dalam analisis data terapan, ada peningkatan ketersediaan kumpulan data yang berisi sejumlah besar peubah. Model linier yang menyertakan lengkap peubah respon seringkali memiliki kinerja prediksi yang buruk karena cenderung memiliki varians yang besar. Selain itu, model besar pada umumnya sulit untuk diinterpretasikan. Dalam banyak kasus, jumlah peubah bahkan lebih besar dari jumlah pengamatan. Metode tradisional seperti *least squares* kemudian tidak dapat lagi diterapkan karena kekurangan rank dari matriks desain.

Untuk meningkatkan akurasi prediksi dan sebagai solusi untuk masalah komputasi dengan data berdimensi tinggi, istilah penalti pada koefisien regresi dapat ditambahkan ke fungsi objektif. Pendekatan ini mengecilkan koefisien dan mengurangi varians dengan harga bias yang meningkat. Tibshirani (1996) memperkenalkan *Least Absolute Shrinkage and Selection Operator* (LASSO), di mana fungsi penalti adalah norma L_1 . Misalkan $\mathbf{y} = (y_1, \dots, y_n)'$ adalah respon dan $\mathbf{X} = (x_{ij})_{1 \leq i \leq n, 1 \leq j \leq p}$ matriks peubah bebas, di mana n menyatakan banyaknya pengamatan dan p menyatakan banyaknya peubah. Selain itu, misalkan x_1, \dots, x_n adalah pengamatan dimensi- p , yaitu baris \mathbf{X} . Kita asumsikan model regresi standar

$$\mathbf{y} = \mathbf{X}^T \boldsymbol{\beta} + e_i \quad (2.46)$$

di mana parameter regresi adalah $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)'$, dan *error* e_i memiliki nilai harapan nol. Dengan parameter penalti λ , penduga LASSO dari $\boldsymbol{\beta}$ adalah

$$\hat{\boldsymbol{\beta}}_{LASSO} = \underset{\boldsymbol{\beta}}{\operatorname{argmin}} \sum_{i=1}^n (\mathbf{y} - \mathbf{X}^T \boldsymbol{\beta})^2 + n\lambda \sum_{j=1}^p |\beta_j| \quad (2.47)$$

LASSO sering digunakan dalam praktik karena penalti L_1 memungkinkan untuk mengecilkan beberapa koefisien hingga tepat nol, yaitu, untuk menghasilkan

estimasi model *sparse* (model di mana hanya sebagian kecil parameter yang tidak nol) yang sangat dapat diinterpretasikan. Selain itu, algoritma cepat untuk menghitung LASSO tersedia melalui kerangka *Least Angle Regression* (LARS). Namun, LASSO tidak *robust* untuk *outlier*. Oleh karena itu, diperlukan alternatif *robust*.

Penduga model *sparse* adalah topik yang sangat penting dalam analisis data modern karena semakin tersedianya kumpulan data dengan sejumlah besar peubah. Metode *Least Trimmed Squares* (LTS) digunakan untuk menemukan himpunan bagian (*subset*) dari h pengamatan yang kuadrat terkecilnya menghasilkan jumlah kuadrat galat terkecil. Ukuran *subset* h dapat dilihat sebagai tebakan awal dari jumlah observasi yang baik dalam data. Sementara LTS bersifat sangat *robust*, akan tetapi tidak menghasilkan estimasi model *sparse*. Selanjutnya, jika $h < p$, estimator LTS tidak dapat dihitung. Versi *sparse* dari LTS diperoleh dengan menambahkan penalti L_1 dengan parameter penalti λ ke persamaan (2.3), sehingga penduga untuk LASSO LTS adalah sebagai berikut,

$$\hat{\beta}_{\text{LASSO-LTS}} = \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^h e_{(i)}^2 + h\lambda \sum_{j=1}^p |\beta_j| \quad (2.48)$$

LASSO LTS memiliki *breakdown point* yang tinggi dan tahan terhadap *outlier* regresi berganda, termasuk *leverage point*. Selain bersifat sangat *robust*, dan mirip dengan penduga LASSO, LASSO LTS memiliki kelebihan

1. Meningkatkan kinerja prediksi melalui pengurangan varians jika ukuran sampel relatif kecil terhadap dimensi
2. Memastikan interpretasi yang lebih tinggi karena pemilihan model simultan
3. Menghindari masalah komputasi metode regresi *robust* tradisional dalam kasus data berdimensi tinggi.

2.8.1 Breakdown Point

Ukuran yang paling populer untuk kekekaran (*robustness*) estimator adalah *replacement finite-sample breakdown* (FBP). Misalkan $\mathbf{Z} = (\mathbf{X}, \mathbf{y})$ menyatakan sampel. Untuk estimator regresi $\hat{\beta}$, *breakdown point* didefinisikan sebagai

$$e^*(\hat{\beta}; \mathbf{Z}) = \min \left\{ \frac{m}{n} : \sup \|\hat{\beta}(\bar{\mathbf{Z}})\|_2 = \infty \right\} \quad (2.49)$$

di mana $\bar{\mathbf{Z}}$ adalah data rusak yang diperoleh dari \mathbf{Z} dengan mengganti m dari n titik data asli dengan nilai *arbitrary*.

Misalkan $\rho(x)$ adalah *loss function* cembung dan simetris dengan $\rho(0) = 0$ dan $\rho(x) > 0$ untuk $x \neq 0$, dan didefinisikan $\boldsymbol{\rho}(x) = (\rho(x_1), \dots, \rho(x_n))^T$. Dengan ukuran himpunan bagian $h \leq n$, maka estimator regresi sebagai berikut:

$$\hat{\boldsymbol{\beta}} = \underset{\boldsymbol{\beta}}{\operatorname{argmin}} \sum_{i=1}^h \boldsymbol{\rho}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})_{i:n} + h\lambda \sum_{j=1}^p |\beta_j| \quad (2.50)$$

di mana $(\boldsymbol{\rho}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}))_{i:n} \leq \dots \leq (\boldsymbol{\rho}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}))_{n:n}$ adalah *regression loss* yang diurutkan. Kemudian *breakdown point* penduga $\hat{\boldsymbol{\beta}}$ diberikan sebagai berikut

$$e^*(\hat{\boldsymbol{\beta}}; \mathbf{Z}) = \frac{n - h + 1}{n} \quad (2.51)$$

Breakdown point sama untuk setiap *loss function* p yang memenuhi asumsi. Secara khusus, *breakdown point* untuk estimator LASSO LTS $\hat{\boldsymbol{\beta}}_{\text{LASSO-LTS}}$ dengan ukuran *subset* $h \leq n$, di mana $\rho(x) = x^2$, adalah $(n - h + 1)/n$. Semakin kecil nilai h , semakin tinggi *breakdown point*-nya. Dengan mengambil h yang cukup kecil, bahkan dimungkinkan untuk memiliki *breakdown point* yang lebih besar dari 50%. Namun, sementara ini secara matematis mungkin, tidak disarankan untuk menggunakan $h < n/2$ karena statistik yang kekar (*robust*) bertujuan untuk model yang sesuai dengan sebagian besar data. Sebagai gantinya, disarankan untuk mengambil nilai h yang sama dengan sebagian kecil dari ukuran sampel, dengan $\alpha = 0.75$, sehingga perkiraan akhir didasarkan pada jumlah pengamatan yang cukup besar.

2.8.2 Algoritma LASSO LTS

Untuk parameter penalti tetap λ berdasarkan persamaan (2.16), maka ditentukan *objective function*

$$Q(H, \boldsymbol{\beta}) = \sum_{i \in H} (y_i - \mathbf{X}'_i \boldsymbol{\beta})^2 + h\lambda \sum_{j=1}^p |\beta_j| \quad (2.52)$$

yang merupakan jumlah kuadrat galat L_1 terpenalti berdasarkan subsample $H \subseteq \{1, \dots, n\}$ dengan $|H| = h$. Dengan

$$\hat{\beta}_H = \underset{\beta}{\operatorname{argmin}}(H, \beta) \quad (2.53)$$

penduga LASSO LTS diberikan oleh $\hat{\beta}_{H_{opt}}$, dimana

$$H_{opt} = \underset{H \subseteq \{1, \dots, n\}: |H|=h}{\operatorname{argmin}} (H, \hat{\beta}_H) \quad (2.54)$$

Oleh karena itu *Sparse* LTS digunakan untuk menemukan himpunan bagian dari $h \leq n$ pengamatan yang *Least Absolute Shrinkage and Selection Operator* (LASSO) *fit*-nya menghasilkan jumlah kuadrat galat terkecil. Untuk menemukan *subset* optimal ini, digunakan analog dari algoritma FAST-LTS yang dikembangkan oleh Rousseeuw dan Van Driessen (2006).

Algoritma ini didasarkan pada *concentration steps* atau *C-steps*. *C-steps* pada iterasi k terdiri dari menghitung solusi LASSO berdasarkan *subset* H_k , dengan $|H_k| = h$, dan membuat *subset* berikutnya H_{k+1} dari pengamatan yang sesuai dengan h kuadrat galat terkecil. Misalkan H_k menyatakan subsampel tertentu yang diturunkan pada iterasi k dan misalkan $\hat{\beta}_{H_k}$ adalah koefisien dari LASSO *fit* yang sesuai. Setelah menghitung kuadrat galat $e_k^2 = (e_{k,1}^2, \dots, e_{k,n}^2)'$ dengan $e_{k,i}^2 = (y_i - \mathbf{X}_i' \hat{\beta}_{H_k})^2$, subsampel H_{k+1} untuk iterasi $k+1$ didefinisikan sebagai himpunan indeks yang sesuai dengan h kuadrat galat terkecil. Dalam istilah matematika, ini dapat ditulis sebagai

$$H_{k+1} = \left\{ i \in \{1, \dots, n\} : e_{k,i}^2 \in \left\{ (e_k^2)_{j:n} : j = 1, \dots, h \right\} \right\} \quad (2.56)$$

Dimana $(e_k^2)_{1:n} \leq \dots \leq (e_k^2)_{n:n}$ menyatakan kuadrat galat yang diurutkan. Misalkan $\hat{\beta}_{H_{k+1}}$ menyatakan koefisien LASSO *fit* berdasarkan H_{k+1} . Maka

$$Q(H_{k+1}, \hat{\beta}_{H_{k+1}}) \leq Q(H_{k+1}, \hat{\beta}_{H_k}) \leq Q(H_k, \hat{\beta}_{H_k}) \quad (2.57)$$

dimana pertidaksamaan pertama mengikuti definisi $\hat{\beta}_{H_{k+1}}$, dan pertidaksamaan kedua dari definisi H_k . Dari persamaan (2.58) didapat bahwa *C-steps* menghasilkan penurunan fungsi objektif LTS yang *sparse*, dan bahwa urutan dari *C-steps* menghasilkan konvergensi ke minimum lokal dalam sejumlah langkah yang terbatas. Untuk meningkatkan peluang mencapai minimum global, jumlah s yang cukup besar dari subsampel awal H_0 harus digunakan, masing-masing digunakan sebagai titik awal untuk urutan *C-steps*. Daripada memilih h titik data secara acak, setiap *subset* awal H_0 dengan ukuran h disusun dari sebuah *elemental subset* berukuran 3. Ambil

tiga pengamatan dari data secara acak, misalkan, x_{i_1} , x_{i_2} , dan x_{i_3} . Maka LASSO *fit* untuk *elemental subset* ini sebagai berikut,

$$\hat{\beta}_{\{i_1, i_2, i_3\}} = \underset{\beta}{\operatorname{argmin}} Q(\{i_1, i_2, i_3\}, \beta) \quad (2.58)$$

dan *subset* awal H_0 kemudian diberikan oleh indeks h pengamatan dengan kuadrat galat terkecil berdasarkan pada persamaan (2.58). Algoritma FAST-LTS *nonsparse* menggunakan himpunan bagian elemen berukuran p , karena setiap regresi OLS memerlukan setidaknya observasi sebanyak dimensi p . Hal ini akan membuat algoritma tidak berlaku jika $p > n$. Untungnya LASSO sudah didefinisikan dengan benar untuk sampel ukuran 3, bahkan untuk nilai p yang besar. Selain itu, dari sudut pandang kekekaran (*robustness*), hanya menggunakan tiga pengamatan yang optimal, karena memastikan probabilitas tertinggi untuk tidak menyertakan *outlier* dalam himpunan *elemental*. Penting untuk dicatat bahwa *elemental subset* berukuran 3 hanya digunakan untuk membangun himpunan bagian awal berukuran h untuk algoritma *C-steps*. Semua *C-steps* dilakukan pada himpunan bagian dari ukuran h .

2.8.3 Reweighted LASSO LTS Estimator

Misalkan α menyatakan proporsi pengamatan dari sampel penuh yang akan dipertahankan di setiap subsampel, dan $h = (n + 1)\alpha$. Kemudian $(1 - \alpha)$ dapat diartikan sebagai tebakan awal dari proporsi *outlier* dalam data. Tebakan awal ini biasanya agak konservatif untuk memastikan bahwa *outlier* tidak mempengaruhi hasil, dan karena itu dapat mengakibatkan hilangnya efisiensi statistik. Untuk meningkatkan efisiensi, langkah *reweighted* yang menurunkan bobot *outlier* yang terdeteksi oleh estimator LASSO LTS dapat dilakukan.

Di bawah model *error* normal, pengamatan dengan galat terstandarisasi yang lebih besar dari kuantil tertentu dari distribusi normal baku dapat dinyatakan sebagai *outlier*. Karena estimasi LASSO LTS seperti LASSO adalah bias, kita perlu memusatkan galatnya. Estimasi natural untuk pusat galat adalah

$$\hat{\mu}_{raw} = \frac{1}{h} \sum_{i \in H_{opt}} e_i \quad (2.59)$$

di mana $e_i = y_i - x_i^T \hat{\beta}_{LASSO-LTS}$ dan H_{opt} adalah himpunan bagian optimal dari (2.54). Kemudian estimasi skala galat yang terkait dengan penduga *raw* LASSO LTS diberikan oleh

$$\hat{\sigma}_{raw} = k_\alpha \sqrt{\frac{1}{h} \sum_{i \in H_{opt}} (e_c^2)_{i:n}} \quad (2.60)$$

dengan kuadrat galat terpusat $e_c^2 = ((e_1 - \hat{\mu}_{raw})^2, \dots, (e_n - \hat{\mu}_{raw})^2)^T$, dan

$$k_\alpha = \left(\frac{1}{\alpha} \int_{-\Phi^{-1}((\alpha+1)/2)}^{\Phi^{-1}((\alpha+1)/2)} u^2 d\Phi(u) \right)^{-1/2} \quad (2.61)$$

k_α merupakan sebuah faktor untuk memastikan bahwa $\hat{\sigma}_{raw}$ adalah estimasi yang konsisten dari standar deviasi pada model normal. Formulasi ini memungkinkan untuk menentukan bobot biner

$$w_i \begin{cases} 1, & \text{jika } |(e_i - \hat{\mu}_{raw})/\hat{\sigma}_{raw}| \leq \Phi^{-1}(1 - \delta), \\ 0, & \text{jika } |(e_i - \hat{\mu}_{raw})/\hat{\sigma}_{raw}| > \Phi^{-1}(1 - \delta), \end{cases} \quad i = 1, \dots, n \quad (2.62)$$

Reweighted estimator LASSO LTS diberikan dengan LASSO *fit* terbobot

$$\hat{\beta}_{reweighted} = \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^n w_i (y_i - x_i^T \beta)^2 + \lambda n_w \sum_{j=1}^p |\beta_j| \quad (2.63)$$

dengan $n_w = \sum_{i=1}^n w_i$ jumlah bobot. Dengan pilihan bobot yang diberikan pada (2.62), *Reweighted* penduga LASSO LTS adalah LASSO *fit* berdasarkan pengamatan tidak ditandai sebagai pencilan. Tentu saja, skema pembobotan lain dapat dipertimbangkan. Menggunakan estimasi pusat galat

$$\hat{\mu}_{reweighted} = \frac{1}{n_w} \sum_{i=1}^n w_i (y_i - x_i^T \hat{\beta}_{reweighted}) \quad (2.64)$$

estimasi skala galat dari *Reweighted* penduga LASSO LTS diberikan oleh

$$\hat{\sigma}_{reweighted} = k_{\alpha_w} \sqrt{\frac{1}{n_w} \sum_{i=1}^n w_i (y_i - x_i^T \hat{\beta}_{reweighted})^2} \quad (2.65)$$

di mana k_{α_w} adalah faktor konsistensi dari persamaan (2.28) dengan $\alpha_w = n_w/n$.

Perhatikan bahwa langkah pembobotan ulang ini secara konseptual berbeda dari *adaptive* LASSO oleh Zou (2006). Sementara *adaptive* LASSO mendapatkan penalti secara individu pada prediktor dari koefisien penduga awal, *Reweighted*

penduga LASSO LTS bertujuan untuk memasukkan semua pengamatan *nonoutlying* ke dalam pencocokan model.

2.8.4 Pemilihan Parameter Penalti

Dalam analisis data praktis, nilai parameter penalti λ yang sesuai tidak diketahui sebelumnya. Disarankan untuk memilih λ dengan mengoptimalkan *Bayes Information Criterion* (BIC), atau perkiraan kinerja prediksi melalui *cross validation* (CV). Dalam penelitian ini digunakan BIC karena membutuhkan lebih sedikit upaya komputasi. BIC dari model yang diestimasi dengan parameter penyusutan λ diberikan oleh

$$\text{BIC}(\lambda) = \log(\hat{\sigma}) + df(\lambda) \frac{\log(n)}{n} \quad (2.66)$$

di mana $\hat{\sigma}$ menunjukkan perkiraan skala galat yang sesuai, (2.60) atau (2.65), dan $df(\lambda)$ adalah derajat kebebasan model. Derajat kebebasan diberikan oleh jumlah parameter estimasi bukan nol di $\hat{\beta}$ (Zou *et al.*, 2007).

Sebagai alternatif BIC, validasi silang dapat digunakan. Untuk mencegah *outlier* mempengaruhi pemilihan λ , prediksi *robust loss function* harus digunakan. Pemilihan natural adalah *root trimmed mean squared prediction error* (RTMSPE) dengan proporsi pemangkasan yang sama seperti untuk menghitung LASSO LTS. Dalam validasi silang *k-fold*, data dibagi secara acak dalam k blok dengan ukuran yang kira-kira sama. Setiap blok dikeluarkan sekali untuk kecocokan model, dan blok yang dikeluarkan digunakan sebagai data uji. Dengan cara ini, dan untuk nilai λ tertentu, sebuah prediksi diperoleh untuk setiap pengamatan dalam sampel. Dinyatakan vektor kuadrat *error* dugaan $\mathbf{e}^2 = (e_1^2, \dots, e_n^2)$. Maka

$$\text{RTMSPE}(\lambda) = \sqrt{\frac{1}{h} \sum_{i \in H_{opt}} (\mathbf{e}^2)_{i:n}} \quad (2.67)$$

Untuk mengurangi variabilitas, RTMSE dirata-ratakan pada lebih dari 500 pemisahan acak data yang berbeda.

Parameter λ yang dipilih kemudian meminimalkan $\text{BIC}(\lambda)$ atau $\text{RTMSPE}(\lambda)$ pada *grid* nilai di interval $[0, \hat{\lambda}_0]$. Diambil *grid* dengan langkah ukuran $0.025\hat{\lambda}_0$, di mana $\hat{\lambda}_0$ adalah perkiraan parameter penyusutan λ_0 yang akan menyusutkan semua

parameter menjadi nol. Jika $p > n$, 0 tentu saja dikeluarkan dari *grid*. Untuk solusi LASSO diperoleh

$$\hat{\lambda}_0 = \frac{2}{n} \max_{j \in \{1, \dots, p\}} \text{Cor}(\mathbf{y}, \mathbf{x}_j) \quad (2.68)$$

$\text{Cor}(\mathbf{y}, \mathbf{x}_j)$ merupakan korelasi Pearson antara \mathbf{y} dan kolom ke- j dari matriks \mathbf{X} . Untuk LASSO LTS, diperlukan estimasi *robust* $\hat{\lambda}_0$. Disarankan untuk mengganti korelasi Pearson pada persamaan (2.68) dengan korelasi kuat berdasarkan *bivariate winorization* data (Khan *et al.*, 2007).

2.9 Tuberkulosis

Di seluruh dunia, sekitar 10 juta orang jatuh sakit tuberkulosis (TB) setiap tahun. TB adalah salah satu dari 10 penyebab kematian teratas, dan penyebab utama dari agen infeksi tunggal (*Mycobacterium tuberculosis*), peringkat di atas HIV/AIDS. Penyakit ini dapat menyerang siapa saja di mana saja, tetapi kebanyakan orang yang mengembangkan TB (sekitar 90%) adalah orang dewasa, rasio pria dan wanita adalah 2:1, dan tingkat kasus di tingkat nasional bervariasi dari kurang dari 50 hingga lebih dari 5000 per 1 juta penduduk. per tahun. Hampir 90% kasus setiap tahun berada di 30 negara dengan beban TB tinggi. Secara global, diperkirakan 1.7 miliar orang terinfeksi *Mycobacterium tuberculosis* dan dengan demikian berisiko terkena penyakit tersebut (WHO, 2018).

Dengan diagnosis dan pengobatan yang tepat waktu serta antibiotik, kebanyakan orang yang terkena TB dapat disembuhkan dan penularan selanjutnya dibatasi. Jumlah kasus yang terjadi setiap tahun (dan dengan demikian jumlah kematian terkait TB) juga dapat diturunkan dengan mengurangi prevalensi faktor risiko terkait kesehatan untuk TB, misalnya merokok, diabetes dan infeksi HIV, memberikan pengobatan pencegahan kepada orang dengan infeksi TB laten, dan tindakan pada determinan infeksi dan penyakit TB yang lebih luas misalnya kemiskinan, kualitas perumahan dan kekurangan gizi (WHO, 2018).

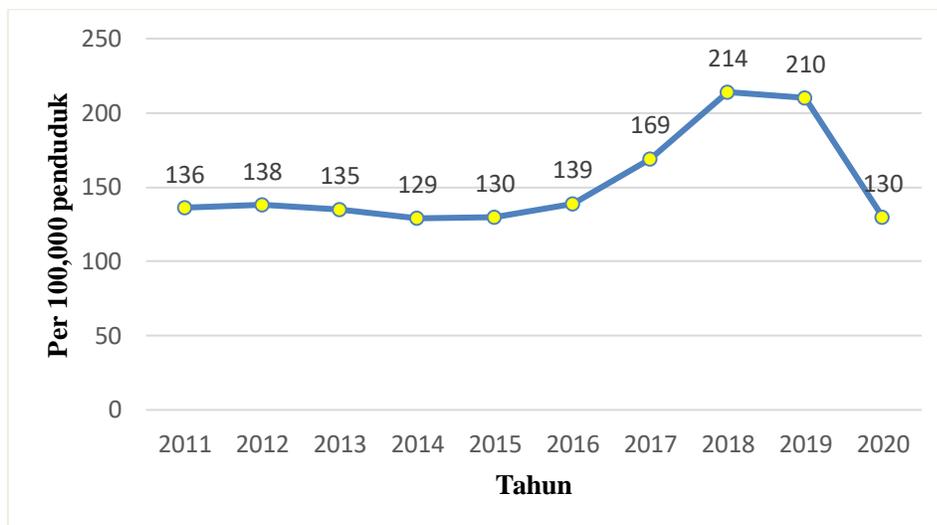
TB saat ini masih merupakan masalah kesehatan masyarakat baik di Indonesia maupun internasional sehingga menjadi salah satu tujuan pembangunan Kesehatan berkelanjutan (SDGs). TB merupakan salah satu dari 10 penyebab utama kematian di seluruh dunia (WHO, 2018).

Indonesia berada pada peringkat ke-3 dengan penderita TB tertinggi di Dunia setelah India dan China. Meskipun terjadi penurunan kasus baru TB, tetapi tidak cukup cepat untuk mencapai target Strategi END TB tahun 2020, yaitu pengurangan kasus TB sebesar 20% antara tahun 2015 – 2020. Pada tahun 2015 – 2019 penurunan kumulatif kasus TB hanya sebesar 9% (WHO, 2018).

Begitu juga dengan kematian akibat TB, jumlah kematian pada tahun 2018 sebesar 1,2 juta. Secara global kematian akibat TB per tahun menurun secara global, tetapi tidak mencapai target Strategi END TB tahun 2020 sebesar 35% antara tahun 2015 – 2020. Jumlah kematian kumulatif antara tahun 2015 – 2019 sebesar 14%, yaitu kurang dari setengah dari target yang ditentukan (WHO, 2018). Beban penyakit yang disebabkan oleh TB dapat diukur salah satunya dengan angka notifikasi semua kasus tuberkulosis atau *Case Notification Rate* (CNR).

2.9.1 *Case Notification Rate* (CNR)

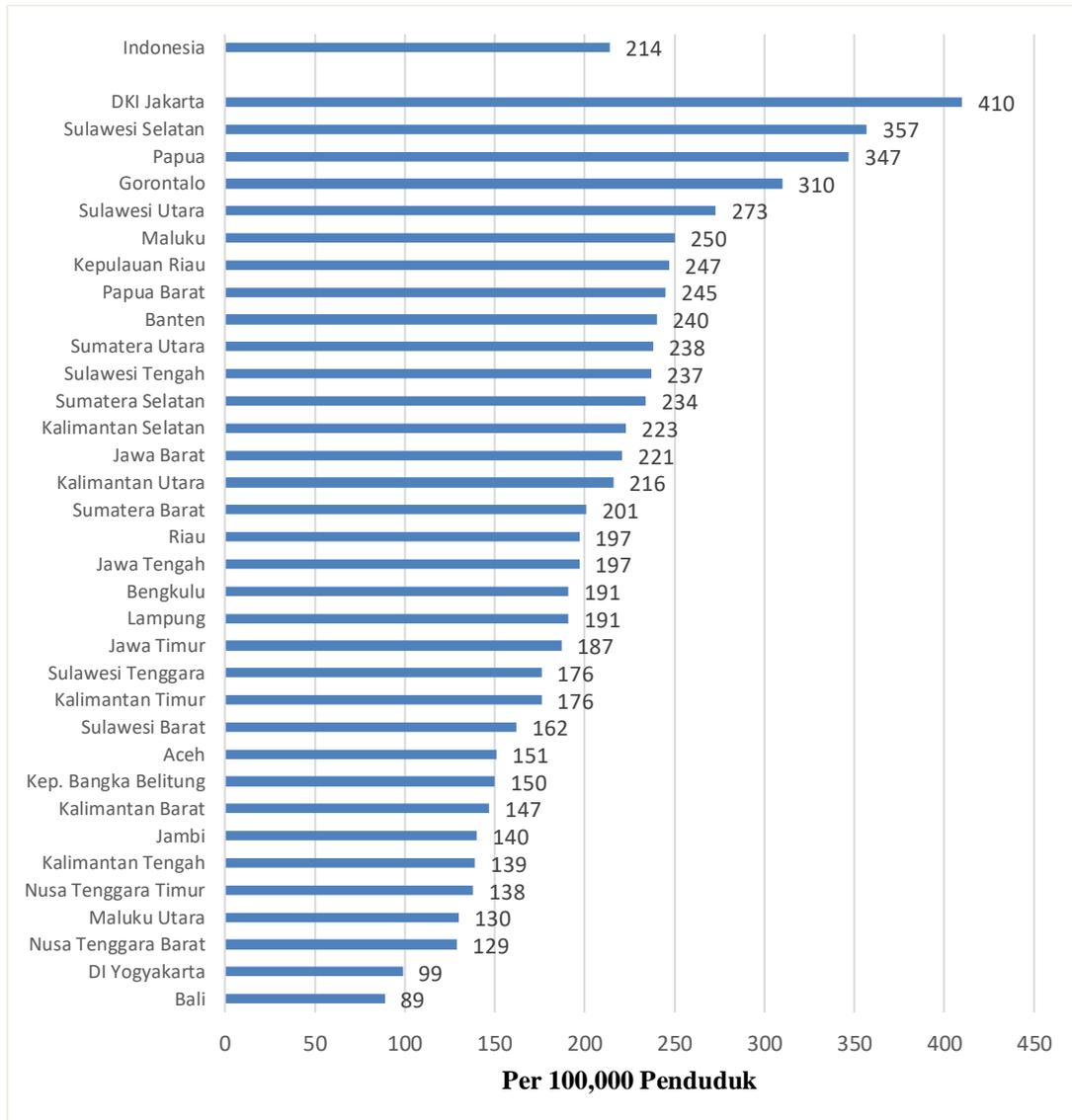
Angka notifikasi Semua Kasus Tuberkulosis atau *Case Notification Rate* (CNR) adalah jumlah semua kasus TB yang diobati dan dilaporkan di antara 100,000 penduduk yang ada di suatu wilayah tertentu. Angka ini apabila dikumpulkan serial, akan menggambarkan kecenderungan (*tren*) meningkat atau menurunnya penemuan kasus dari tahun ke tahun di suatu wilayah (Kemenkes RI, 2021).



Gambar 1. Angka Notifikasi Semua Kasus TB Per 100,000 Penduduk Indonesia Tahun 2011-2020

Sumber: Kemenkes RI, 2021

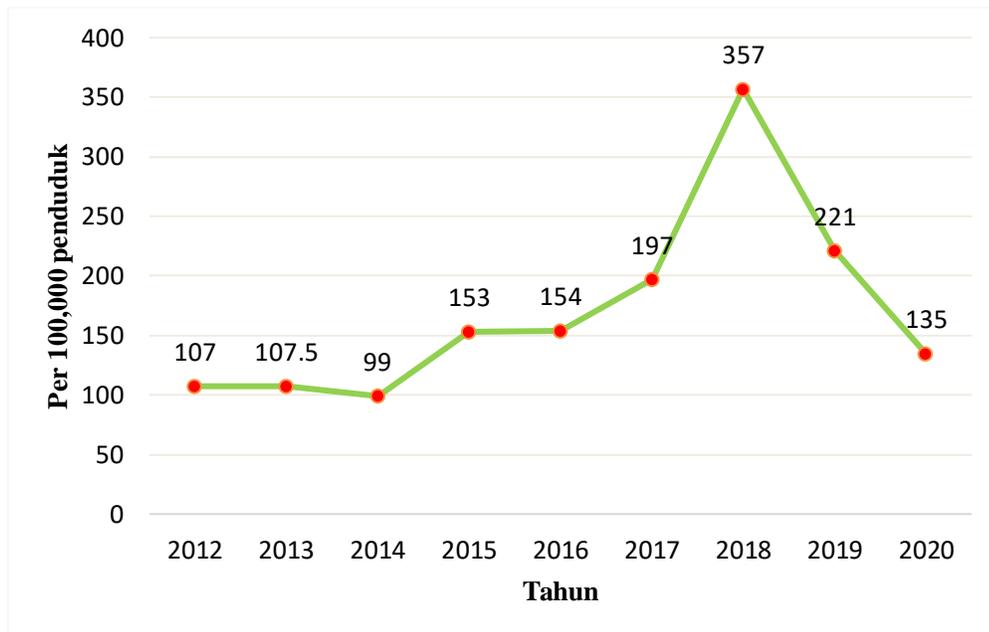
Gambar 1. menunjukkan angka notifikasi semua kasus TB per 100,000 penduduk dari tahun 2010-2020 yang secara nasional memperlihatkan kecenderungan peningkatan CNR sampai tahun 2018 dan menurun pada tahun 2019 dan 2020.



Gambar 2. Angka Notifikasi Semua Kasus TB Per 100,000 Penduduk Menurut Provinsi Tahun 2018

Sumber: Kemenkes RI, 2018

Gambar 2. menunjukkan cakupan semua kasus TB (per 100,000 penduduk) menurut Provinsi pada tahun 2018. Salah satu Provinsi dengan CNR yang tertinggi di Indonesia adalah Provinsi Sulawesi Selatan yaitu sebesar 357 per 100,000 penduduk (Kemenkes RI, 2018).

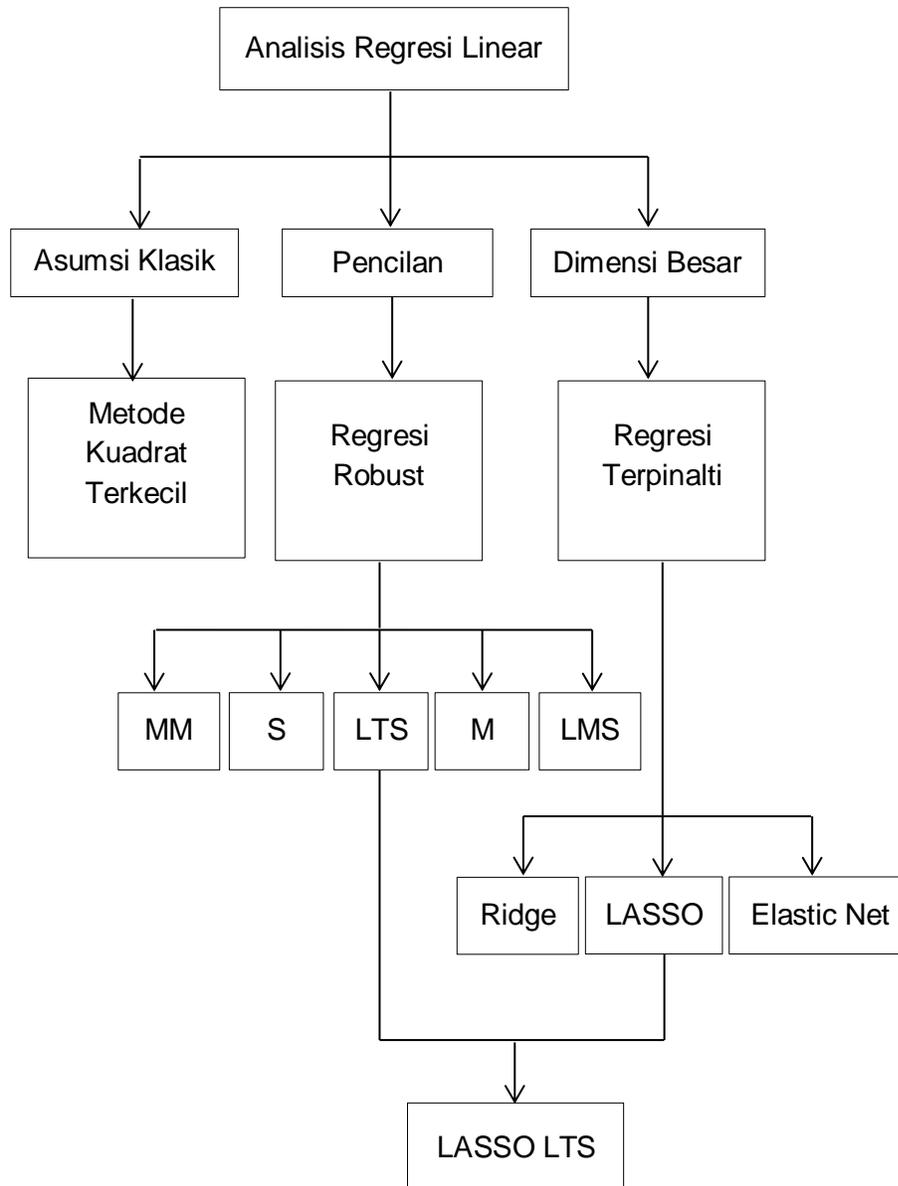


Gambar 3. Angka Notifikasi Semua Kasus TB Per 100.000 Penduduk Provinsi Sulawesi Selatan Tahun 2012-2020

Sumber: Kemenkes RI, 2020

Gambar 3. menunjukkan angka notifikasi semua kasus TB per 100,000 penduduk dari tahun 2012-2020 untuk Provinsi Sulawesi Selatan yang juga memperlihatkan kecenderungan peningkatan CNR sampai tahun 2018 dan menurun pada tahun 2019 dan 2020. Namun, penurunan jumlah kasus tuberkulosis masih tetap harus diwaspadai karena adanya disparitas (kesenjangan) penyebaran penyakit tuberkulosis antar daerah di Sulawesi Selatan, misalnya kabupaten Sidrap pada tahun 2018 memiliki jumlah kasus tuberkulosis sebanyak 493 kasus, namun pada tahun 2020 mengalami kenaikan jumlah kasus sebesar 591 kasus. Hal tersebut diduga karena datanya mengandung pencilan (*outliers*). Oleh sebab itu, perlu dilakukan penelitian lebih lanjut mengenai faktor-faktor yang mempengaruhi penyebaran penyakit tuberkulosis di Provinsi Sulawesi Selatan.

2.10 Kerangka Konseptual



Gambar 4. Diagram Kerangka Konseptual Metode Penelitian