

Condition Number Analysis of Logistic Regression, and its Implications for First-Order Solution Methods

Robert M. Freund (MIT)

joint with Paul Grigas (Berkeley) and Rahul Mazumder (MIT)

CMU, Tepper School of Business, May 2019

How can optimization inform statistics (and machine learning)?

This talk is based on our paper:

Condition Number Analysis of Logistic Regression, and its Implications for First-Order Solution Methods

A “cousin” paper of ours:

A New Perspective on Boosting in Linear Regression via Subgradient Optimization and Relatives

Outline

- Optimization primer: three basic first-order methods for convex optimization
- Logistic regression perspectives: statistics “vs.” machine learning
- A pair of condition numbers for the logistic regression problem:
 - when the sample data is **non-separable**:
 - a condition number for the degree of non-separability of the dataset
 - informing the convergence guarantees of Greedy Coordinate Descent and Stochastic Gradient Descent (SGD)
 - guarantees on reaching linear convergence (thanks to Bach)
 - when the sample data is **separable**:
 - a condition number for the degree of separability of the dataset
 - informing convergence guarantees to deliver an approximate maximum margin classifier

Optimization

Three Basic First-Order Methods for Convex Optimization:

- Greedy Coordinate Descent method
“go in the best coordinate direction”
- Gradient Descent method
“go in the direction of the negative of the gradient”
- Stochastic Gradient Descent (SGD) method
“go in the direction of the negative of the stochastic estimate of the gradient”

Convex Optimization

The problem of interest is:

$$F^* := \min_x F(x) \\ \text{s.t. } x \in \mathbb{R}^p$$

where $F(\cdot)$ is differentiable and convex:

$$F(\lambda x + (1 - \lambda)y) \leq \lambda F(x) + (1 - \lambda)F(y) \quad \text{for all } x, y, \text{ and all } \lambda \in [0, 1]$$

Let $\|x\|$ denote the given norm on the variables $x \in \mathbb{R}^p$

Norms and Dual Norms

Let $\|x\|$ be the given norm on the variables $x \in \mathbb{R}^p$

The dual norm is $\|s\|_* := \max_x \{s^T x : \|x\| \leq 1\}$

Some common norms and their dual norms:

Name	Norm	Definition	Dual Norm
ℓ_2 -norm	$\ x\ _2$	$\ x\ _2 = \sqrt{\sum_{j=1}^p x_j ^2}$	$\ s\ _* = \ s\ _2$
ℓ_1 -norm	$\ x\ _1$	$\ x\ _1 = \sum_{j=1}^p x_j $	$\ s\ _* = \ s\ _\infty$
ℓ_∞ -norm	$\ x\ _\infty$	$\ x\ _\infty = \max\{ x_1 , \dots, x_p \}$	$\ s\ _* = \ s\ _1$

Lipschitz constant for the Gradient

$$F^* := \min_x F(x)$$

$$\text{s.t. } x \in \mathbb{R}^p$$

We say that $\nabla F(\cdot)$ is Lipschitz with parameter L_F if:

$$\|\nabla F(x) - \nabla F(y)\|_* \leq L_F \|x - y\| \quad \text{for all } x, y \in \mathbb{R}^p$$

 $\|\cdot\|_*$ is the dual norm

Matrix Operator Norm

Let M be a linear operator (matrix) $M : \mathbb{R}^p \rightarrow \mathbb{R}^n$ with norm $\|x\|_a$ on \mathbb{R}^p and norm $\|v\|_b$ on \mathbb{R}^n

The operator norm of M is given by:

$$\|M\|_{a,b} := \max_{x \neq 0} \frac{\|Mx\|_b}{\|x\|_a}$$

Greedy Coordinate Descent Method: “go in the best coordinate direction”

$$F^* := \min_x F(x)$$

$$\text{s.t. } x \in \mathbb{R}^p$$

Greedy Coordinate Descent

Initialize at $x^0 \in \mathbb{R}^p$, $k \leftarrow 0$

At iteration k :

- 1 Compute gradient $\nabla F(x^k)$
- 2 Compute
 - $j_k \in \arg \max_{j \in \{1, \dots, p\}} \{|\nabla F(x^k)_j|\}$ and
 - $d^k \leftarrow \text{sgn}(\nabla F(x^k)_{j_k}) e_{j_k}$
- 3 Choose step-size α_k
- 4 Set $x^{k+1} \leftarrow x^k - \alpha_k d^k$

Greedy Coordinate Descent \equiv Steepest Descent in the ℓ_1 -Norm

$$F^* := \min_x F(x) \\ \text{s.t. } x \in \mathbb{R}^p$$

Steepest Descent method in the ℓ_1 -norm

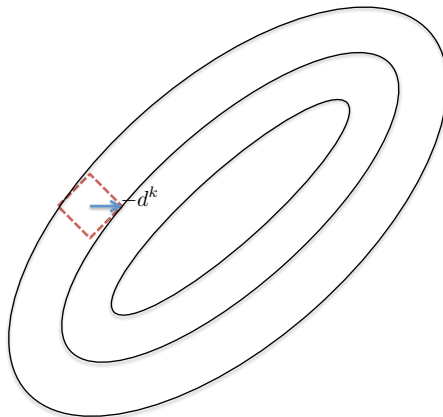
Initialize at $x^0 \in \mathbb{R}^p$, $k \leftarrow 0$

At iteration k :

- ① Compute gradient $\nabla F(x^k)$
- ② Compute direction: $d^k \leftarrow \arg \max_{\|d\|_1 \leq 1} \{\nabla F(x^k)^T d\}$
- ③ Choose step-size α_k
- ④ Set $x^{k+1} \leftarrow x^k - \alpha_k d^k$

Greedy Coordinate Descent \equiv Steepest Descent in the ℓ_1 -Norm, cont.

$$d^k \leftarrow \arg \max_{\|d\|_1 \leq 1} \{\nabla F(x^k)^T d\}$$



Metrics for Evaluating Greedy Coordinate Descent

$$F^* := \min_x F(x) \\ \text{s.t. } x \in \mathbb{R}^p$$

Assume $F(\cdot)$ is convex and $\nabla F(\cdot)$ is Lipschitz with parameter L_F :

$$\|\nabla F(x) - \nabla F(y)\|_\infty \leq L_F \|x - y\|_1 \quad \text{for all } x, y \in \mathbb{R}^p$$

Two sets of interest:

$\mathcal{S}_0 := \{x \in \mathbb{R}^p : F(x) \leq F(x^0)\}$ is the level set of the initial point x^0

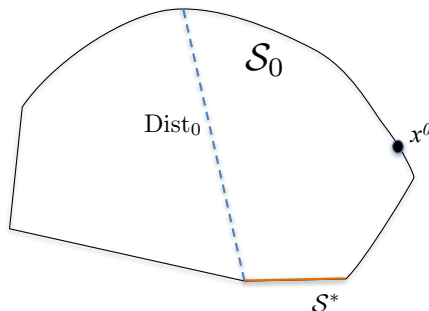
$\mathcal{S}^* := \{x \in \mathbb{R}^p : F(x) = F^*\}$ is the set of optimal solutions

Metrics for Evaluating Greedy Coordinate Descent, cont.

$\mathcal{S}_0 := \{x \in \mathbb{R}^p : F(x) \leq F(x^0)\}$ is the level set of the initial point x^0

$\mathcal{S}^* := \{x \in \mathbb{R}^p : F(x) = F^*\}$ is the set of optimal solutions

$$\text{Dist}_0 := \max_{x \in \mathcal{S}_0} \min_{x^* \in \mathcal{S}^*} \|x - x^*\|_1$$



(In high-dimensional machine learning problems, \mathcal{S}^* can be very big)

Computational Guarantees for Greedy Coordinate Descent

$$\text{Dist}_0 := \max_{x \in \mathcal{S}_0} \min_{x^* \in \mathcal{S}^*} \|x - x^*\|_1$$

Theorem: Objective Function Value Convergence (essentially [Beck and Tetrushvili 2014], [Nesterov 2003])

If the step-sizes are chosen using the rule:

$$\alpha_k = \frac{\|\nabla F(x^k)\|_\infty}{L_F} \quad \text{for all } k \geq 0,$$

then for each $k \geq 0$ the following inequality holds:

$$F(x^k) - F^* \leq \frac{2L_F(\text{Dist}_0)^2}{\hat{K}^0 + k} < \frac{2L_F(\text{Dist}_0)^2}{k}$$

where $\hat{K}^0 := \frac{2L_F(\text{Dist}_0)^2}{F(x^0) - F^*}.$

Computational Guarantees for GCD, cont.

Theorem: Gradient Norm Convergence

For any step-size sequence $\{\alpha_k\}$ and for each $k \geq 0$, it holds that:

$$\min_{i \in \{0, \dots, k\}} \|\nabla F(x^i)\|_\infty \leq \frac{F(x^0) - F^* + \frac{L_F}{2} \sum_{i=0}^k \alpha_i^2}{\sum_{i=0}^k \alpha_i}.$$

If the step-sizes are chosen using the rule:

$$\alpha_k = \frac{\|\nabla F(x^k)\|_\infty}{L_F} \quad \text{for all } k \geq 0,$$

then for each $k \geq 0$ the following inequality holds:

$$\min_{i \in \{0, \dots, k\}} \|\nabla F(x^i)\|_\infty \leq \sqrt{\frac{2L_F(F(x^0) - F^*)}{k+1}}.$$

Computational Guarantees for GCD, cont.

Theorem: Iterate Shrinkage

For any step-size sequence $\{\alpha_k\}$, it holds for each $k \geq 0$ that:

$$\|x^k - x^0\|_1 \leq \sum_{i=0}^{k-1} \alpha_i .$$

If the step-sizes are chosen using the rule:

$$\alpha_k = \frac{\|\nabla F(x^k)\|_\infty}{L_F} \quad \text{for all } k \geq 0 ,$$

then for each $k \geq 0$ it holds that:

$$\|x^k - x^0\|_1 \leq \sqrt{k} \sqrt{\frac{2(F(x^0) - F^*)}{L_F}} .$$

Gradient Descent $\equiv \ell_2$ -Steepest Descent

The problem of interest is:

$$F^* := \min_x F(x) \\ \text{s.t. } x \in \mathbb{R}^p$$

where $F(x)$ is convex and differentiable.

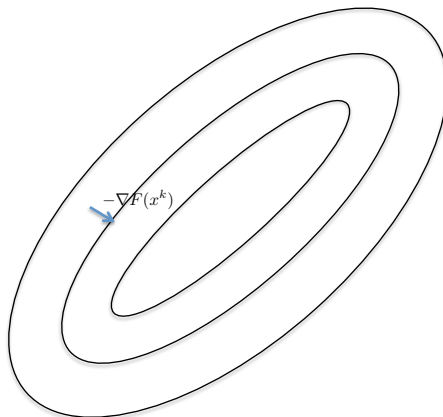
Gradient Descent method for minimizing $f(x)$

Initialize at $x^0 \in \mathbb{R}^p$, $k \leftarrow 0$

At iteration k :

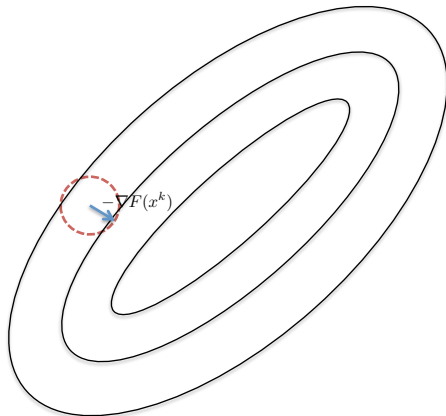
- 1 Compute gradient $\nabla F(x^k)$
- 2 Choose step-size $\hat{\alpha}_k$
- 3 Set $x^{k+1} \leftarrow x^k - \hat{\alpha}_k \nabla F(x^k)$

Gradient Descent $\equiv \ell_2$ -Steepest Descent, cont.



Gradient Descent $\equiv \ell_2$ -Steepest Descent, cont.

$$\frac{\nabla F(x^k)}{\|\nabla F(x^k)\|_2} \in \arg \max_{\|d\|_2 \leq 1} \{\nabla F(x^k)^T d\}$$



Stochastic Gradient Descent (SGD) Method

The problem of interest is:

$$F^* := \min_x F(x) \\ \text{s.t. } x \in \mathbb{R}^p$$

Let $\tilde{\nabla} f(x)$ be a stochastic estimate of the gradient $\nabla F(x)$ at each x

Stochastic Gradient Descent method for minimizing $F(x)$

Initialize at $x^0 \in \mathbb{R}^p$, $k \leftarrow 0$

At iteration k :

- 1 Compute stochastic gradient $\tilde{\nabla} F(x^k)$
- 2 Choose step-size α_k
- 3 Set $x^{k+1} \leftarrow x^k - \alpha_k \tilde{\nabla} F(x^k)$

Stochastic Gradient Descent (SGD) Method, cont.

The canonical setting for SGD is minimizing a large sum (or average) of losses:

$$\begin{aligned} F^* &:= \min_x F(x) := \frac{1}{n} \sum_{j=1}^n F_j(x) \\ \text{s.t. } &x \in \mathbb{R}^p \end{aligned}$$

where $n \gg 0$ and $\tilde{\nabla} F(x)$ is computed as follows:

- 1 Choose $j \sim \{1, \dots, n\}$ uniformly and independently
- 2 $\tilde{\nabla} F(x) \leftarrow \nabla F_j(x)$

Then the stochastic gradient is unbiased: $\mathbb{E}[\tilde{\nabla} F(x)|x] = \nabla F(x)$

Computational Guarantees for SGD [Bertsekas, others]

Assume that

(i) the stochastic gradient is unbiased, namely

$$\mathbb{E}[\tilde{\nabla} F(x)|x] = \nabla F(x) \text{ for any } x, \text{ and}$$

(ii) $F(\cdot)$ is G -stochastically smooth: there exists G such that:

$$\mathbb{E}[\|\tilde{\nabla} F(x)\|_2^2 | x] \leq G^2 \text{ for any } x$$

Theorem: Expected Convergence of Stochastic Gradient Descent

If the step-sizes are constant:

$$\alpha_k = \bar{\alpha} \text{ for all } k \geq 0,$$

then for each $k \geq 0$ the following inequality holds:

$$\mathbb{E}[F(\bar{x}^k)] - F^* \leq \frac{\bar{\alpha} G^2}{2} + \frac{\|x^0 - x^*\|_2^2}{2\bar{\alpha}(k+1)},$$

where $\bar{x}^k := \frac{1}{k+1} \sum_{i=0}^k x^i$.

Logistic Regression

Logistic Regression

- statistics perspective
- machine learning perspective

Logistic Regression: Statistics Perspective

Logistic Regression: Statistics Perspective

Logistic Regression Statistics Perspective

Example: Predicting Parole Violation

Predict $P(\text{violate parole})$ based on age, gender, time served, offense class, multiple convictions, NYC, etc.

	Violator	Male	Age	TimeServed	Class	Multiple	InCity
1	0	1	49.4	3.15	D	0	1
2	1	1	26.0	5.95	D	1	0
3	0	1	24.9	2.25	D	1	0
4	0	1	52.1	29.22	A	0	0
5	0	1	35.9	12.78	A	1	1
6	0	1	25.9	1.18	C	1	1
7	0	1	19.0	0.54	D	0	0
8	0	1	43.2	1.07	C	0	1
9	0	1	31.6	1.17	E	0	0
10	0	1	40.7	4.64	B	1	1
11	0	1	53.9	21.61	A	0	1
12	0	1	28.5	3.23	D	1	0
13	0	1	36.1	3.71	D	0	1
14	0	1	48.8	1.17	D	0	0
15	0	1	37.6	4.62	C	0	0
16	0	1	42.5	1.75	D	0	1
...
6098	0	1	55.0	0.72	E	0	0
6099	0	1	49.6	29.88	A	0	1
6100	0	1	22.4	2.85	D	0	1
6101	0	1	44.8	1.76	D	1	0
6102	0	0	45.3	1.03	E	0	0

Logistic Regression for Prediction

$Y \in \{-1, 1\}$ is a Bernoulli random variable:

$$P(Y = 1) = p$$

$$P(Y = -1) = 1 - p$$

$x = (x_1, \dots, x_p) \in \mathbb{R}^p$ is the vector of independent variables

$P(Y = 1)$ depends on the values of the independent variables x_1, \dots, x_p

Logistic regression model is:

$$P(Y = 1 \mid x) = \frac{1}{1 + e^{-\beta^T x}}$$

Logistic Regression for Prediction, continued

Logistic regression model is:

$$P(Y = 1 \mid x) = \frac{1}{1 + e^{-\beta^T x}}$$

Data records are (x_i, y_i) , $i = 1, \dots, n$

	Violator	Male	Age	TimeServed	Class	Multiple	InCity
1	0	1	49.4	3.15	D	0	1
2	1	1	26.0	5.95	D	1	0
3	0	1	24.9	2.25	D	1	0
4	0	1	52.1	29.22	A	0	0
5	0	1	35.9	12.78	A	1	1
6	0	1	25.9	1.18	C	1	1
7	0	1	19.0	0.54	D	0	0
8	0	1	43.2	1.07	C	0	1
9	0	1	31.6	1.17	E	0	0
10	0	1	40.7	4.64	B	1	1
11	0	1	53.9	21.61	A	0	1
12	0	1	28.5	3.23	D	1	0
13	0	1	36.1	3.71	D	0	1
14	0	1	48.8	1.17	D	0	0
15	0	1	37.6	4.62	C	0	0
16	0	1	42.5	1.75	D	0	1
...
6098	0	1	55.0	0.72	E	0	0
6099	0	1	49.6	29.88	A	0	1
6100	0	1	22.4	2.85	D	0	1
6101	0	1	44.8	1.76	D	1	0
6102	0	0	45.3	1.03	E	0	0

Let us construct an estimate of β based on the data (x_i, y_i) , $i = 1, \dots, n$

Logistic Regression: Maximum Likelihood Estimation

$$\begin{aligned}
 & \max_{\beta} \left(\prod_{y_i=1} \frac{1}{1 + e^{-\beta^T x_i}} \right) \left(\prod_{y_i=-1} \left(1 - \frac{1}{1 + e^{-\beta^T x_i}} \right) \right) \\
 &= \max_{\beta} \left(\prod_{i=1}^n \frac{1}{1 + e^{-y_i \beta^T x_i}} \right) \\
 &\equiv \min_{\beta} \frac{1}{n} \sum_{i=1}^n \ln \left(1 + e^{-y_i \beta^T x_i} \right) =: L_n(\beta)
 \end{aligned}$$

Logistic Regression Optimization Problem

Logistic regression optimization problem is:

$$\begin{aligned} L_n^* &:= \min_{\beta} L_n(\beta) := \frac{1}{n} \sum_{i=1}^n \ln(1 + \exp(-y_i \beta^T x_i)) \\ \text{s.t. } &\beta \in \mathbb{R}^p \end{aligned}$$

If $y_i = +1$, we ideally want $\beta^T x_i \gg 0$

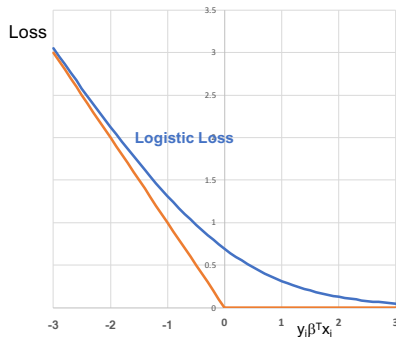
If $y_i = -1$, we ideally want $\beta^T x_i \ll 0$

Therefore we ideally want β for which $y_i \beta^T x_i \gg 0$ for very many i

Logistic Regression Optimization Problem, continued

Logistic regression optimization problem is:

$$\begin{aligned} L_n^* &:= \min_{\beta} L_n(\beta) := \frac{1}{n} \sum_{i=1}^n \ln(1 + \exp(-y_i \beta^T x_i)) \\ \text{s.t. } &\beta \in \mathbb{R}^p \end{aligned}$$



Each logistic loss term is a 1-smoothing of $\max\{0, -y_i \beta^T x_i\}$

Properties of the Logistic Loss Function

$$L_n^* := \min_{\beta} L_n(\beta) := \frac{1}{n} \sum_{i=1}^n \ln(1 + \exp(-y_i \beta^T x_i))$$

s.t. $\beta \in \mathbb{R}^p$

Denote $\mathbf{X} := \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}$

Proposition: Lipschitz constant of the gradient of $L_n(\beta)$

$\nabla L_n(\cdot)$ is $L = \frac{1}{4n} \|\mathbf{X}\|_{1,2}^2$ -Lipschitz:

$$\|\nabla L_n(\beta) - \nabla L_n(\beta')\|_{\infty} \leq \frac{1}{4n} \|\mathbf{X}\|_{1,2}^2 \|\beta - \beta'\|_1$$

where $\|\mathbf{X}\|_{1,2} := \max_{\|\beta\|_1 \leq 1} \|\mathbf{X}\beta\|_2$

Properties of the Logistic Loss Function, continued

$$L_n^* := \min_{\beta} L_n(\beta) := \frac{1}{n} \sum_{i=1}^n \ln(1 + \exp(-y_i \beta^T x_i))$$

s.t. $\beta \in \mathbb{R}^p$

- $L_n(\cdot)$ is convex
- $L_n^* \geq 0$
- If $L_n^* = 0$, then the optimum is not attained (something is “wrong” or “very wrong”)
- We will see later that “very wrong” is actually very good....
- For $\beta^0 := 0$ it holds that $L_n(\beta^0) = \ln(2)$

Logistic Regression: Machine Learning Perspective

Logistic Regression: Machine Learning Perspective

Logistic Regression Machine Learning Perspective

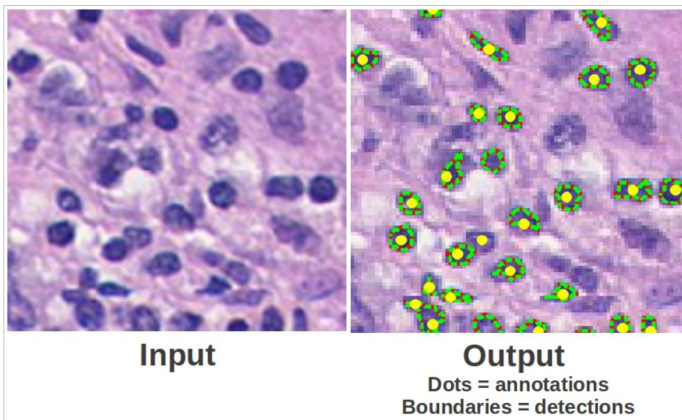
Example: Gender Classification

Classify (predict) gender based on image



Another Example: Cancer/noncancerous cells

Classify (predict) cancer/noncancer cells based on image



Another Example: Voters-nonvoters

Classify (predict) voters vs. nonvoters based on election data

Observation	Age	Income (\$K/Year)	Number of Children	Gender (Female = 1)	Voting Status
1	54	81	1	1	Voted
2	29	68	0	0	Did not Vote
3	42	106	2	1	Voted
4	74	55	1	1	Voted
5	65	75	0	0	Voted
6	35	102	3	0	Voted
7	23	29	0	0	Did not Vote
8	40	36	1	1	Did not Vote
9	24	69	2	0	Did not Vote
10	66	82	0	1	Voted
11	61	94	2	0	Voted
12	36	60	0	1	Did not Vote
13	53	25	1	0	Did not Vote
14	22	72	1	0	Did not Vote
15	41	133	2	1	Voted
16	47	82	3	1	Voted
17	28	37	1	0	Did not Vote
18	43	50	2	0	Did not Vote
19	37	135	4	1	Voted
20	20	58	0	1	Did not Vote
21	62	60	1	0	Voted
22	48	29	2	0	Did not Vote
...
623,151	22	41	0	1	Did not Vote

Binary Classification

Data: $(x_i, y_i) \in \mathbb{R}^p \times \{-1, 1\}$, $i = 1, \dots, n$

- $x = (x_1, \dots, x_p) \in \mathbb{R}^p$ is the vector of features (indep. variables)
- $y \in \{-1, 1\}$ is the set of possible responses/labels

Task: predict y based on the linear function $\beta^T x$

- $\beta \in \mathbb{R}^p$ are the model coefficients

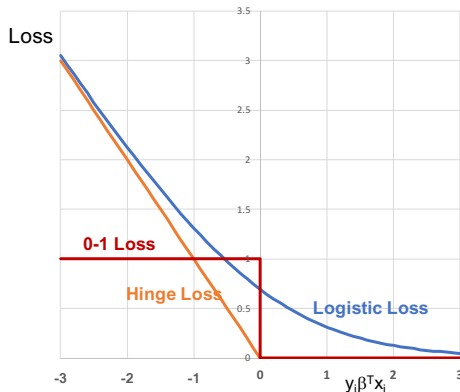
Loss function: $\ell(y, \beta^T x)$ represents the loss incurred when the truth is y but our classification/prediction is based on $\beta^T x$

Empirical Loss Minimization Problem: $\min_{\beta} \frac{1}{n} \sum_{i=1}^n \ell(y_i, \beta^T x_i)$

Loss Functions for Binary Classification

Some common loss functions used for binary classification

- **0-1 loss:** $\ell(y, \beta^T x) := \mathbf{1}(y\beta^T x < 0)$
- **Hinge loss:** $\ell(y, \beta^T x) := \max\{0, -y\beta^T x\}$
- **Logistic loss:** $\ell(y, \beta^T x) := \ln(1 + \exp(-y\beta^T x))$



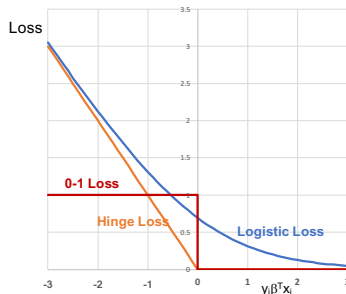
Advantages of Logistic Loss Function

Why use the logistic loss function for classification?

- Computational advantages: convex, smooth
- Fits previous statistical model of conditional probability:

$$P(Y = y \mid x) = \frac{1}{1 + \exp(-y\beta^T x)}$$

- Makes sense when the data is non-separable
- Robust to misspecification of class labels



Logistic Regression Problem of Interest, continued

Alternate versions of optimization problem add regularization and/or sparsification:

$$L_n^* := \min_{\beta} L_n(\beta) := \frac{1}{n} \sum_{i=1}^n \ln(1 + \exp(-y_i \beta^T x_i)) + \lambda \|\beta\|_p$$

s.t. $\beta \in \mathbb{R}^p$

$$\|\beta\|_0 \leq k$$

Overall aspirations:

- Good predictive performance on new (out of sample) observations
- Models that are more interpretable (e.g., sparse)

Computational Experiment: Greedy Coordinate Descent (GCD)

Consider Greedy Coordinate Descent (GCD) for Logistic Regression

Greedy Coordinate Descent for Logistic Regression

Greedy Coordinate Descent for Logistic Regression

Initialize at $\beta^0 \leftarrow 0, k \leftarrow 0$

At iteration $k \geq 0$:

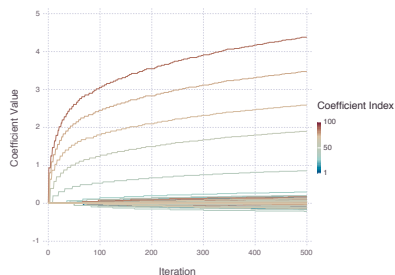
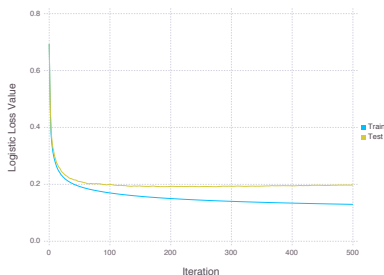
- 1 Compute $\nabla L_n(\beta^k)$
- 2 Compute $j_k \in \arg \max_{j \in \{1, \dots, p\}} |\nabla L_n(\beta^k)_j|$
- 3 Set $\beta^{k+1} \leftarrow \beta^k - \alpha_k \text{sgn}(\nabla L_n(\beta^k)_{j_k}) e_{j_k}$

Why use Greedy Coordinate Descent for Logistic Regression?

- Scalable and effective when $n, p \gg 0$ and maybe $p > n$
- GCD performs variable selection
- GCD imparts implicit regularization
- Just one tuning parameter (number of iterations)

Implicit Regularization and Variable Selection Properties

Artificial example: $n = 1000, p = 100$, true model has 5 non-zeros



Compare with explicit regularization schemes (ℓ_1, ℓ_2 , etc.)

How do GCD and SGD Inform Logistic Regression?

Some questions:

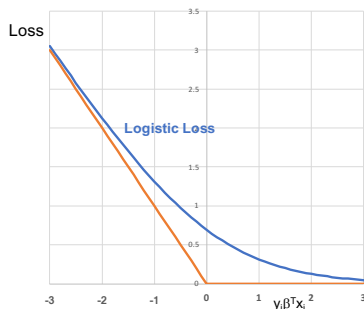
- How do the computational guarantees for Greedy Coordinate Descent and Stochastic Gradient Descent specialize to the case of Logistic Regression?
- Can we say anything further about the convergence properties of these methods in the special case of Logistic Regression?
- What role does problem structure/conditioning play in these guarantees?

Elementary Properties of the Logistic Loss Function

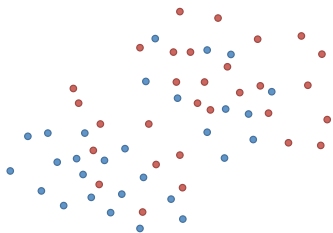
$$L_n^* := \min_{\beta} L_n(\beta) := \frac{1}{n} \sum_{i=1}^n \ln(1 + \exp(-y_i \beta^T x_i))$$

Recall that logistic regression “ideally” seeks β for which $y_i x_i^T \beta \gg 0$ for all i :

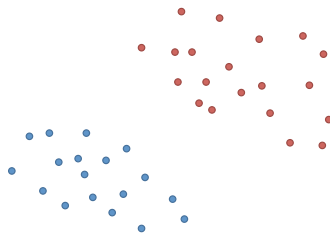
- $y_i = +1 \Rightarrow x_i^T \beta \gg 0$
- $y_i = -1 \Rightarrow x_i^T \beta \ll 0$



Geometry of the Data: Separable and Non-Separable Data

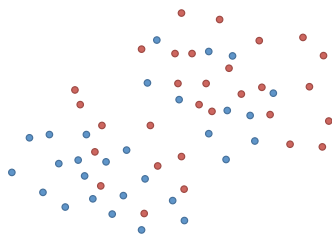


(a) Data is Non-Separable

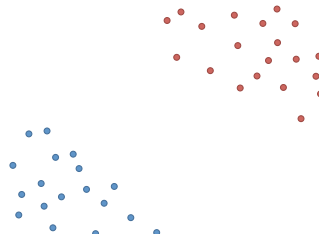


(b) Data is Separable

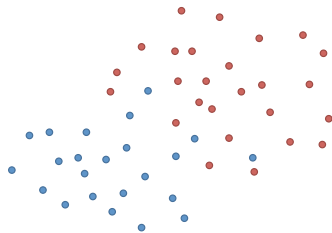
Very/Mild Separable/Non-Separable Data



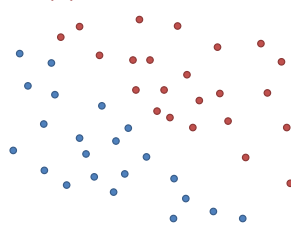
(a) Data is Very Non-Separable



(b) Data is Very Separable



(c) Data is Mildly Non-Separable



(d) Data is Mildly Separable

Separable and Non-Separable Data

Separable Data

The data is separable if there exists $\bar{\beta}$ for which

$$y_i \cdot (\bar{\beta})^T x_i > 0 \quad \text{for all } i = 1, \dots, n$$

Non-Separable Data

The data is non-separable if it is not separable, namely, every β satisfies

$$y_i \cdot (\beta)^T x_i \leq 0 \quad \text{for at least one } i \in \{1, \dots, n\}$$

Separable Data and Non-Attainment of Optimum

$$L_n^* := \min_{\beta} L_n(\beta) := \frac{1}{n} \sum_{i=1}^n \ln(1 + \exp(-y_i \beta^T x_i))$$

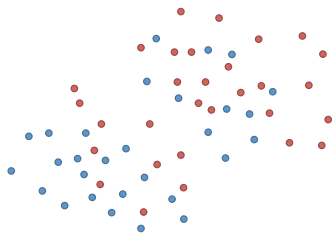
The data is separable if there exists $\bar{\beta}$ for which

$$y_i \cdot (\bar{\beta})^T x_i > 0 \quad \text{for all } i = 1, \dots, n$$

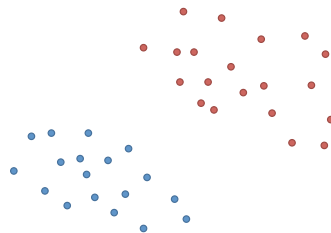
If $\bar{\beta}$ separates the data, then $L_n(\theta \bar{\beta}) \rightarrow 0 (= L_n^*)$ as $\theta \rightarrow +\infty$

Perhaps trying to optimize the logistic loss function is unlikely to be effective at finding a “good” linear classifier

Separable and Non-Separable Data



(a) Data is Non-Separable



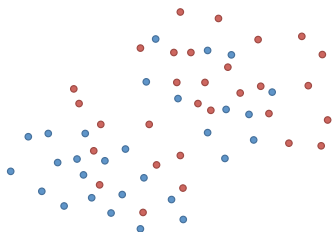
(b) Data is Separable

Results in the Non-Separable Case

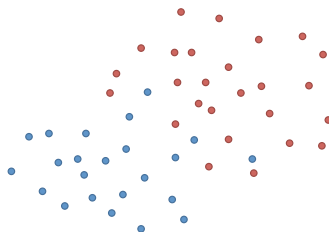
Results in the Non-Separable Case

Non-Separable Data and Problem Behavior/Conditioning

Let us quantify the degree of non-separability of the data.



(a) Very non-separable data



(b) Mildly non-separable data

We will relate this to problem behavior/conditioning....

Non-Separability Condition Number DegNSEP^*

Definition of Non-Separability Condition Number DegNSEP^*

$$\begin{aligned} \text{DegNSEP}^* &:= \min_{\beta \in \mathbb{R}^p} \quad \frac{1}{n} \sum_{i=1}^n [y_i \beta^T x_i]^- \\ &\text{s.t.} \quad \|\beta\|_1 = 1 \end{aligned}$$

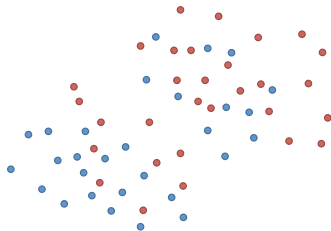
DegNSEP^* is the least average misclassification error (over all normalized classifiers)

$\text{DegNSEP}^* > 0$ if and only if the data is strictly non-separable

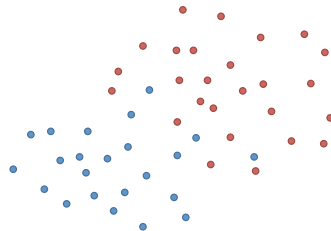
Non-Separability Measure DegNSEP*

$$\text{DegNSEP}^* := \min_{\beta \in \mathbb{R}^p} \quad \frac{1}{n} \sum_{i=1}^n [y_i \beta^T x_i]^-$$

s.t. $\|\beta\|_1 = 1$



(a) DegNSEP* is large



(b) DegNSEP* is small

DegNSEP* and Problem Behavior/Conditioning

$$L_n^* := \min_{\beta} L_n(\beta) := \frac{1}{n} \sum_{i=1}^n \ln(1 + \exp(-y_i \beta^T x_i))$$

$$\begin{aligned} \text{DegNSEP}^* := \min_{\beta \in \mathbb{R}^p} \quad & \frac{1}{n} \sum_{i=1}^n [y_i \beta^T x_i]^- \\ \text{s.t.} \quad & \|\beta\|_1 = 1 \end{aligned}$$

Theorem: Non-Separability and Sizes of Optimal Solutions

Suppose that the data is non-separable and $\text{DegNSEP}^* > 0$. Then

① the logistic regression problem LR attains its (unique) optimum,

② for every optimal solution β^* of LR it holds that

$$\|\beta^*\|_1 \leq \frac{L_n^*}{\text{DegNSEP}^*} \leq \frac{\ln(2)}{\text{DegNSEP}^*}, \text{ and}$$

③ for any β it holds that $\|\beta\|_1 \leq \frac{L_n(\beta)}{\text{DegNSEP}^*}$.

Algorithmic Results for Non-Separable Case

- Computational Guarantees for Greedy Coordinate Descent
- Reaching Linear Convergence in Greedy Coordinate Descent
- Computational Guarantees for Stochastic Gradient Descent

Computational Guarantees for Greedy Coordinate Descent: Non-Separable Case

Theorem: Computational Guarantees for Greedy Coordinate Descent: Non-Separable Case

Consider the GCD applied to the Logistic Regression problem with step-sizes $\alpha_k := \frac{4n\|\nabla L_n(\beta^k)\|_\infty}{\|\mathbf{X}\|_{1,2}^2}$ for all $k \geq 0$, and suppose that the data is non-separable. Then for each $k \geq 0$ it holds that:

- (i) (training error): $L_n(\beta^k) - L_n^* \leq \frac{2(\ln(2))^2 \|\mathbf{X}\|_{1,2}^2}{k \cdot n \cdot \text{DegNSEP}^*}$
- (ii) (gradient norm): $\|\nabla L_n(\beta^k)\|_\infty \leq \frac{\|\mathbf{X}\|_{1,2}^2 \ln(2)}{\sqrt{k \cdot n \cdot \text{DegNSEP}^*}}$
- (iii) (regularization): $\|\beta^k\|_1 \leq \sqrt{k} \left(\frac{1}{\|\mathbf{X}\|_{1,2}} \right) \sqrt{8n(\ln(2) - L_n^*)}$

Computational Guarantees for Greedy Coordinate Descent: Non-Separable Case

Theorem: Computational Guarantees for Greedy Coordinate Descent: Non-Separable Case

Consider the GCD applied to the Logistic Regression problem with step-sizes $\alpha_k := \frac{4n\|\nabla L_n(\beta^k)\|_\infty}{\|\mathbf{X}\|_{1,2}^2}$ for all $k \geq 0$, and suppose that the data is non-separable. Then for each $k \geq 0$ it holds that:

- (i) (training error): $L_n(\beta^k) - L_n^* \leq \frac{2(\ln(2))^2 \|\mathbf{X}\|_{1,2}^2}{k \cdot n \cdot (\text{DegNSEP}^*)^2}$
- (ii) (gradient norm): $\|\nabla L_n(\beta^k)\|_\infty \leq \frac{\|\mathbf{X}\|_{1,2}^2 \ln(2)}{\sqrt{k} \cdot n \cdot \text{DegNSEP}^*}$
- (iii) (regularization): $\|\beta^k\|_1 \leq \sqrt{k} \left(\frac{1}{\|\mathbf{X}\|_{1,2}} \right) \sqrt{8n(\ln(2) - L_n^*)}$

Reaching Linear Convergence

Reaching Linear Convergence using Gradient Descent for Logistic Regression

For logistic regression, does Gradient Descent exhibit linear convergence?

Some Definitions/Notation

Definitions:

- $R := \max_{i \in \{1, \dots, n\}} \|x_i\|_2$ (maximum ℓ_2 norm of the feature vectors)
- $H(\beta^*)$ denotes the Hessian of $L_n(\cdot)$ at an optimal solution β^*
- $\lambda_{\min}(H(\beta^*))$ denotes the smallest eigenvalue of $H(\beta^*)$
- For this section only, let us replace the ℓ_1 -norm and ℓ_∞ -norm by the ℓ_2 -norm

Reaching Linear Convergence of Gradient Descent for Logistic Regression

Theorem: Reaching Linear Convergence of Gradient Descent for Logistic Regression, the “slow part”

Consider Gradient Descent applied to the Logistic Regression problem with step-sizes $\alpha_k := \frac{4n\|\nabla L_n(\beta^k)\|_2}{\|\mathbf{X}\|_{2,2}^2}$ for all $k \geq 0$, and suppose that the data is non-separable. Define the “slow” rate of linear convergence constant:

$$\tau_{\text{slow}} := \left(1 - \frac{2(\text{DegNSEP}^*)\lambda_{\min}(H(\beta^*))n}{(\text{DegNSEP}^* + 2\ln(2)\|\mathbf{X}\|_{2,\infty})\|\mathbf{X}\|_{2,2}^2} \right) < 1 .$$

Then for all $k \geq 0$, it holds that:

(i) (training error): $L_n(\beta^k) - L_n^* \leq (\ln(2) - L_n^*) \cdot (\tau_{\text{slow}})^k$, and

(ii) (coefficient convergence):

$$\|\beta^k - \beta^*\| \leq \left(1 + \frac{2\ln(2)\|\mathbf{X}\|_{2,\infty}}{\text{DegNSEP}^*} \right) \left(\frac{\|\mathbf{X}\|_{2,2}}{\lambda_{\min}(H(\beta^*))} \right) \sqrt{\frac{\ln(2) - L_n^*}{2n}} \cdot (\tau_{\text{slow}})^{k/2} ,$$

where β^* is the unique optimal solution of LR.

Reaching Linear Convergence of Gradient Descent for Logistic Regression, cont.

Theorem: Reaching Linear Convergence of Gradient Descent for Logistic Regression, the “fast part”

Furthermore, define:

$$\check{K} := \left\lceil \frac{16 \ln(2)^2 \|\mathbf{X}\|_{2,2}^4 \|\mathbf{X}\|_{2,\infty}^2}{9n^2 (\text{DegNSEP}^*)^2 \lambda_{\min}(H(\beta^*))^2} \right\rceil ,$$

and the “fast” rate of linear convergence constant:

$$\tau_{\text{fast}} := \left(1 - \frac{\lambda_{\min}(H(\beta^*))n}{\|\mathbf{X}\|_{2,2}^2} \right) < \tau_{\text{slow}} < 1 .$$

Then for all $k \geq \check{K}$, it holds that:

(iii) (training error): $L_n(\beta^k) - L_n^* \leq (L_n(\beta^{\check{K}}) - L_n^*) \cdot (\tau_{\text{fast}})^{k-\check{K}}$, and

(iv) (coefficient convergence):

$$\|\beta^k - \beta^*\| \leq \frac{\|\mathbf{X}\|_{2,2}}{\lambda_{\min}(H(\beta^*))} \sqrt{\frac{2(L_n(\beta^{\check{K}}) - L_n^*)}{n}} \cdot (\tau_{\text{fast}})^{(k-\check{K})/2} .$$

Reaching Linear Convergence of Gradient Descent for Logistic Regression, cont.

Some comments:

- Proof relies on (a slight generalization of) the “generalized self-concordance” property of the logistic loss function due to [Bach 2014]
- Furthermore, we can bound:

$$\lambda_{\min}(H(\beta^*)) \geq \frac{1}{4n} \lambda_{\min}(\mathbf{X}^T \mathbf{X}) \exp \left(-\frac{\ln(2) \|\mathbf{X}\|_{2,\infty}}{\text{DegNSEP}^*} \right)$$

- As compared to results of a similar flavor for other algorithms, here we have an exact characterization of when the linear convergence “kicks in” and also what the rate of linear convergence is guaranteed to be
- Q: Can we exploit this generalized self-concordance property in other ways? (still ongoing ...)

Computational Guarantees for Stochastic Gradient Descent: Non-Separable Case

Theorem: Computational Guarantees for Stochastic Gradient Descent: Non-Separable Case

Consider SGD applied to the Logistic Regression problem with constant step-size

$$\alpha_i := \bar{\alpha} = \frac{n \ln(2)}{\|\mathbf{X}\|_F^2 \sqrt{k+1}}$$

for $i = 0, \dots, k$, and suppose that the data is non-separable. Then it holds that:

$$(\text{training error}): \mathbb{E}[L_n(\bar{\beta}^k)] - L_n^* \leq \frac{\ln(2)}{2\sqrt{k+1}} \left(\frac{\|\mathbf{X}\|_F^2}{n(\text{DegNSEP}^*)^2} + 1 \right)$$

where $\bar{\beta}^k := \frac{1}{k+1} \sum_{i=0}^k \beta^i$.

Computational Guarantees for Stochastic Gradient Descent: Non-Separable Case

Theorem: Computational Guarantees for Stochastic Gradient Descent: Non-Separable Case

Consider SGD applied to the Logistic Regression problem with constant step-size

$$\alpha_i := \bar{\alpha} = \frac{n \ln(2)}{\|\mathbf{X}\|_F^2 \sqrt{k+1}}$$

for $i = 0, \dots, k$, and suppose that the data is non-separable. Then it holds that:

$$(\text{training error}): \mathbb{E}[L_n(\bar{\beta}^k)] - L_n^* \leq \frac{\ln(2)}{2\sqrt{k+1}} \left(\frac{\|\mathbf{X}\|_F^2}{n(\text{DegNSEP}^*)^2} + 1 \right)$$

where $\bar{\beta}^k := \frac{1}{k+1} \sum_{i=0}^k \beta^i$.

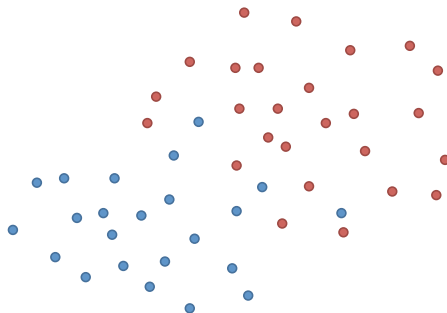
DegNSEP* and “Perturbation to Separability”

$$\begin{aligned} \text{DegNSEP}^* := & \min_{\beta \in \mathbb{R}^p} \quad \frac{1}{n} \sum_{i=1}^n [y_i \beta^T x_i]^- \\ \text{s.t.} \quad & \|\beta\|_1 = 1 \end{aligned}$$

Theorem: DegNSEP* is the “Perturbation to Separability”

$$\begin{aligned} \text{DegNSEP}^* = & \inf_{\Delta x_1, \dots, \Delta x_n} \quad \frac{1}{n} \sum_{i=1}^n \|\Delta x_i\|_\infty \\ \text{s.t.} \quad & (x_i + \Delta x_i, y_i), i = 1, \dots, n \text{ are separable} \end{aligned}$$

Illustration of Perturbation to Separability

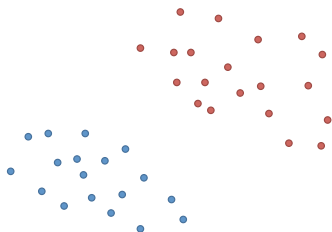


Results in the Separable Case

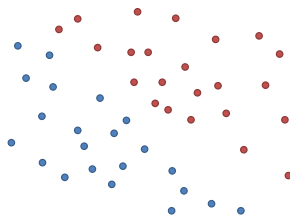
Results in the Separable Case

Separable Data and Problem Behavior/Conditioning

Let us quantify the degree of separability of the data.



(a) Very separable data



(b) Barely separable data

We will relate this to problem behavior/conditioning....

Separability Condition Number DegSEP^*

Definition of Separability Condition Number DegSEP^*

$$\begin{aligned} \text{DegSEP}^* &:= \max_{\beta \in \mathbb{R}^p} \min_{i \in \{1, \dots, n\}} [y_i \beta^T x_i] \\ \text{s.t.} \quad &\|\beta\|_1 \leq 1 \end{aligned}$$

DegSEP^* maximizes the minimal classification value $[y_i \beta^T x_i]$ (over all normalized classifiers)

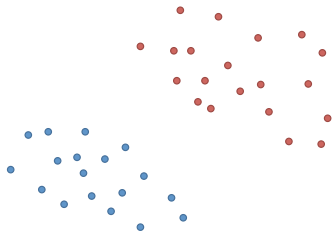
DegSEP^* is simply the “maximum margin” in machine learning parlance

$\text{DegSEP}^* > 0$ if and only if the data is separable

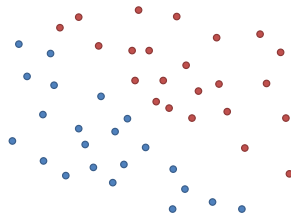
Separability Measure DegSEP*

$$\text{DegSEP}^* := \max_{\beta \in \mathbb{R}^p} \min_{i \in \{1, \dots, n\}} [y_i \beta^T x_i]$$

s.t. $\|\beta\|_1 \leq 1$



(a) DegSEP* is large



(b) DegSEP* is small

DegSEP* and Problem Behavior/Conditioning

$$L_n^* := \min_{\beta} L_n(\beta) := \frac{1}{n} \sum_{i=1}^n \ln(1 + \exp(-y_i \beta^T x_i))$$

$$\begin{aligned} \text{DegSEP}^* &:= \max_{\beta \in \mathbb{R}^p} \min_{i \in \{1, \dots, n\}} [y_i \beta^T x_i] \\ \text{s.t.} \quad &\|\beta\|_1 \leq 1 \end{aligned}$$

Theorem: Separability and Non-Attainment

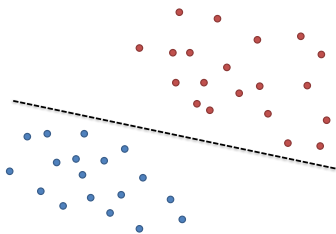
Suppose that the data is separable. Then $\text{DegSEP}^* > 0$, $L_n^* = 0$, and LR does not attain its optimum.

Despite this, it turns out that the Steepest Descent family and also Stochastic Gradient Descent are reasonably effective at finding an approximate margin maximizer as we shall shortly see....

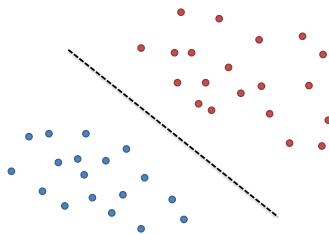
Margin function $\rho(\beta)$

Margin function $\rho(\beta)$

$$\rho(\beta) := \min_{i \in \{1, \dots, n\}} [y_i \beta^T x_i]$$



(a) $\rho(\beta)$ is small



(b) $\rho(\beta)$ is large

Algorithmic Results for the Separable Case

- Computational Guarantees for Gradient Descent
- Computational Guarantees for Stochastic Gradient Descent

Computational Guarantees for Gradient Descent: Separable Case

Theorem: Computational Guarantees for Gradient Descent: Separable Case

Consider Gradient Descent applied to the Logistic Regression problem with step-sizes $\alpha_k := \frac{2\|\nabla L_n(\beta^k)\|_2}{\|\mathbf{X}\|_{2,\infty}^2}$ for all $k \geq 0$, and suppose that the data is separable.

- (i) (margin bound): there exists $i \in \{0, \dots, k\}$ for which the normalized iterate $\bar{\beta}^i := \beta^i / \|\beta^i\|_2$ satisfies

$$\rho(\bar{\beta}^i) \geq \frac{\text{DegSEP}^* \cdot \ln \left(\frac{\text{DegSEP}^*}{n\|\mathbf{X}\|_{2,\infty}} \sqrt{\frac{3(k+1)}{2\ln(2)}} - 1 \right)}{2(\ln(k) + 1)}.$$

- (ii) (shrinkage): $\|\beta^k\|_2 \leq \frac{2\ln(k)}{\text{DegSEP}^*} + \frac{2}{\|\mathbf{X}\|_{2,\infty}}$

- (iii) (gradient bound): $\min_{i \in \{0, \dots, k\}} \|\nabla L_n(\beta^i)\|_2 \leq \|\mathbf{X}\|_{2,\infty} \sqrt{\frac{2\ln(2)}{3(k+1)}}$

Computational Guarantees for Gradient Descent: Separable Case

Theorem: Computational Guarantees for Gradient Descent: Separable Case

Consider Gradient Descent applied to the Logistic Regression problem with step-sizes $\alpha_k := \frac{2\|\nabla L_n(\beta^k)\|_2}{\|\mathbf{X}\|_{2,\infty}^2}$ for all $k \geq 0$, and suppose that the data is separable.

- (i) (margin bound): there exists $i \in \{0, \dots, k\}$ for which the normalized iterate $\bar{\beta}^i := \beta^i / \|\beta^i\|_2$ satisfies

$$\rho(\bar{\beta}^i) \geq \frac{\text{DegSEP}^* \cdot \ln \left(\frac{\text{DegSEP}^*}{n\|\mathbf{X}\|_{2,\infty}} \sqrt{\frac{3(k+1)}{2\ln(2)}} - 1 \right)}{2(\ln(k) + 1)}.$$

- (ii) (shrinkage): $\|\beta^k\|_2 \leq \frac{2\ln(k)}{\text{DegSEP}^*} + \frac{2}{\|\mathbf{X}\|_{2,\infty}}$

- (iii) (gradient bound): $\min_{i \in \{0, \dots, k\}} \|\nabla L_n(\beta^i)\|_2 \leq \|\mathbf{X}\|_{2,\infty} \sqrt{\frac{2\ln(2)}{3(k+1)}}$

Implications for convergence to the margin DegSEP^*

- The previous theorem implies

$$\rho(\bar{\beta}^i) \geq \frac{1}{4} \text{DegSEP}^* \cdot \left(1 - \frac{\ln(2) + \frac{1}{2} \ln \left(\frac{2 \ln(2) n^2 \|\mathbf{X}\|_{2,\infty}^2}{3(\text{DegSEP}^*)^2} \right) + \frac{1}{2}}{\frac{1}{2} \ln(k+1) + \frac{1}{2}} \right)$$

- Except for the factor of $\frac{1}{4}$, this is comparable to Soudry, Hoffer, and Srebro [2017], and improves on Ji and Telgarsky [2018].

Computational Guarantees for Stochastic Gradient Descent: Separable Case

Theorem: Computational Guarantees for Stochastic Gradient Descent: Separable Case

Consider SGD applied to the Logistic Regression problem with step-sizes $\alpha_i := \frac{\ln(2)}{\sqrt{k+1}\|\mathbf{X}\|_{2,\infty}^2}$ for $i = 0, \dots, k$. Choose $\hat{\beta}^k \sim \mathcal{U}[\beta^1, \dots, \beta^k]$. Then:

- (i) (margin bound): For any $\gamma \in (0, 1]$, with probability at least $1 - \gamma$ the normalized iterate $\bar{\beta}^k := \hat{\beta}^k / \|\hat{\beta}^k\|$ satisfies

$$\rho(\bar{\beta}^k) > \frac{\text{DegSEP}^* \cdot \ln \left(\frac{\text{DegSEP}^* \sqrt{\gamma} \sqrt[4]{k+1}}{n\|\mathbf{X}\|_{2,\infty} \sqrt{1.1}} - 1 \right)}{2(\ln(k) + 1)} \quad (1)$$

- (ii) (shrinkage): $\|\hat{\beta}^k\|_2 \leq \frac{2\ln(k)}{\text{DegSEP}^*} + \frac{2}{\|\mathbf{X}\|_{2,\infty}}$ and

- (iii) (expected gradient bound): $\mathbb{E} \left[\|\nabla L_n(\hat{\beta}^k)\|_2^2 \right] < \frac{1.1 \cdot \|\mathbf{X}\|_{2,\infty}^2}{\sqrt{k+1}}$

Implications for convergence to the margin DegSEP^*

- The previous theorem implies

$$\rho(\bar{\beta}^k) \geq \frac{1}{8} \text{DegSEP}^* \cdot \left(1 - \frac{\ln(2) - \frac{1}{2} \ln(\gamma) + \frac{1}{4} \ln \left(\frac{(1.1)^2 n^4 \|\mathbf{X}\|_{2,\infty}^4}{(\text{DegSEP}^*)^4} \right) + \frac{1}{4}}{\frac{1}{4} \ln(k+1)} \right)$$

- Except for the factor of $\frac{1}{8}$, this improves on Nacson, Srebro, Soudry [2018]

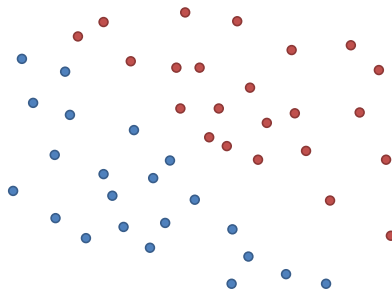
DegSEP* and “Perturbation to Non-Separability”

$$\begin{aligned} \text{DegSEP}^* := & \max_{\beta \in \mathbb{R}^p} \min_{i \in \{1, \dots, n\}} [y_i \beta^T x_i] \\ \text{s.t.} \quad & \|\beta\|_1 \leq 1 \end{aligned}$$

Theorem: DegSEP* is the “Perturbation to Non-Separability”

$$\begin{aligned} \text{DegSEP}^* = & \inf_{\Delta x_1, \dots, \Delta x_n} \max_{i \in \{1, \dots, n\}} \|\Delta x_i\|_\infty \\ \text{s.t.} \quad & (x_i + \Delta x_i, y_i), i = 1, \dots, n \text{ are non-separable} \end{aligned}$$

Illustration of Perturbation to Non-Separability



Other Issues

Some other topics not mentioned (still ongoing):

- Other first-order methods for logistic regression (accelerated gradient descent, other randomized methods, etc.
- High-dimensional regime $p > n$, define DegNSEP_k^* and DegSEP_k^* for restricting β to satisfy $\|\beta\|_0 \leq k$
- Numerical experiments comparing methods
- Other...

Summary

- Some old and new results for Greedy Coordinate Descent, Gradient Descent, and Stochastic Gradient Descent
- Analyzing these methods for Logistic Regression: separable/non-separable cases
- Non-Separable case
 - condition number DegNSEP^*
 - computational guarantees for Greedy Coordinate Descent, Gradient Descent, and Stochastic Gradient Descent, including reaching linear convergence
- Separable case
 - condition number DegSEP^*
 - computational guarantees for Greedy Coordinate Descent, Gradient Descent, and Stochastic Gradient Descent, including computing an approximate maximum margin classifier