

EEP 596: Adv Intro ML || Lecture 4

Dr. Karthik Mohan

Univ. of Washington, Seattle

January 12, 2023

Logistics

- Conceptual 1 is assigned
- Any questions on logistics?

Today's class!

- Recap of Overfitting and Regularization
- Gradient Descent and SGD Algorithm
- Introduction to Classification in ML

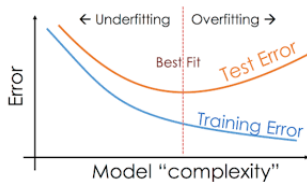
ICE #1 (2 mins)

When is Model A under-fitting as compared to Model B?

Let A_{train} be train error of model A and A_{val} be validation error of model A and the same notation for model B.

- a) $A_{train} < B_{train}$ and $A_{val} > B_{val}$
- b) $A_{train} > B_{train}$ and $B_{val} < A_{val}$
- c) $A_{train} < B_{train}$ and $A_{val} < B_{val}$
- d) $A_{train} > B_{train}$ and $B_{val} > A_{val}$

ICE #1 graphed



Over-fitting and Remedies

Remedy	Name	Benefits
ℓ_2 Reg.	Ridge Regression	No large weights
ℓ_1 Reg.	Lasso	Removes un-important features
$\ell_1 - \ell_2$ Reg.	Elastic Net	Combined benefits
Feature Selection		Reduces d so that $d \ll N$
Increase dataset size	Data Aug.	Increases N so that $N \gg d$

Understanding l_1 and l_2 norms better



$$p = \infty$$



$$p = 2$$

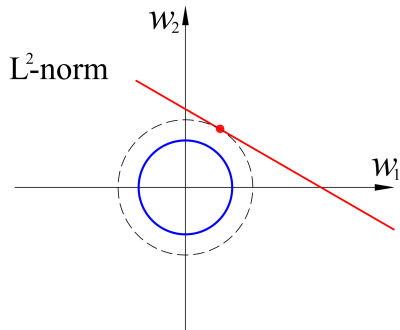
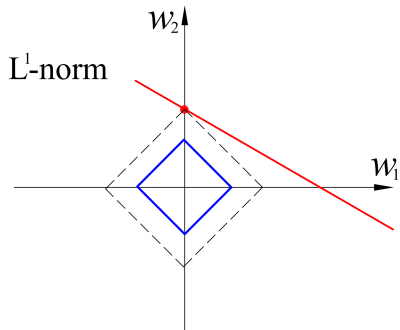


$$p = 1$$

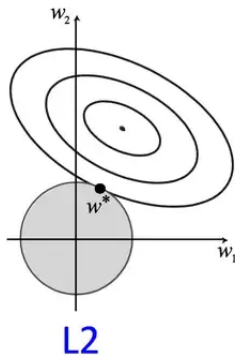
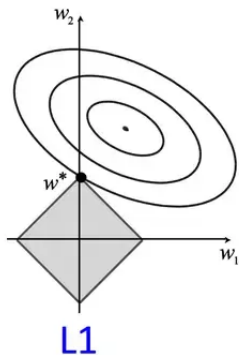


$$0 < p < 1$$

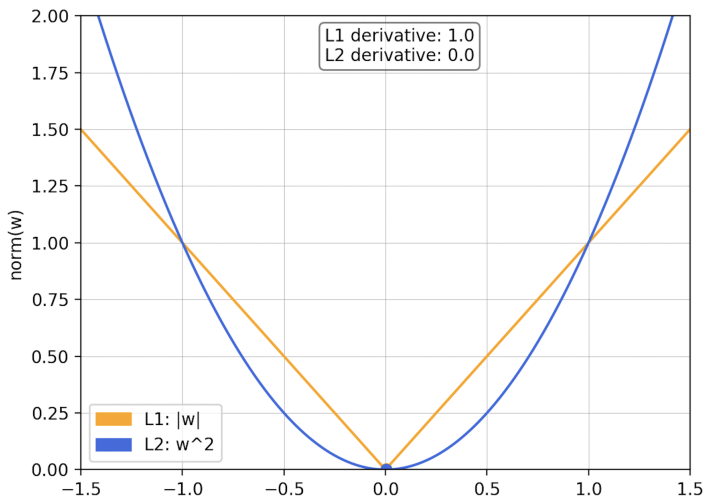
Understanding l_1 and l_2 norms better



Understanding l_1 and l_2 norms better



Understanding l_1 and l_2 norms in one dimension



Over-fitting and Remedies

Remedy	Name	Benefits
ℓ_2 Reg.	Ridge Regression	No large weights
ℓ_1 Reg.	Lasso	Removes un-important features
$\ell_1 - \ell_2$ Reg.	Elastic Net	Combined benefits
Feature Selection		Reduces d so that $d \ll N$
Increase dataset size	Data Aug.	Increases N so that $N \gg d$

ICE #2

Manhattan and Euclidean Distance

Every **norm** of a vector (or a matrix) gives rise to a **distance metrics**. Norm is a measure of magnitude of a vector (or matrix) while distance metric is a measure of well, distance between two vectors. Consider for instance the distance between Seattle and Bellevue. If you drew a straight-line between the two cities, that would be the **Euclidean distance**. However, if you start in downtown seattle, and take SR-520, that is equivalent to the ℓ_1 distance or **Manhattan distance**. Compute the Euclidean and Manhattan distance between two vectors, $x = [1, 2, 3], y = [2, 4, -1]$. The distances are closest to:

- 1 7 and 4
- 2 4 and 7
- 3 7 and 5
- 4 5 and 7

Understanding regularization better

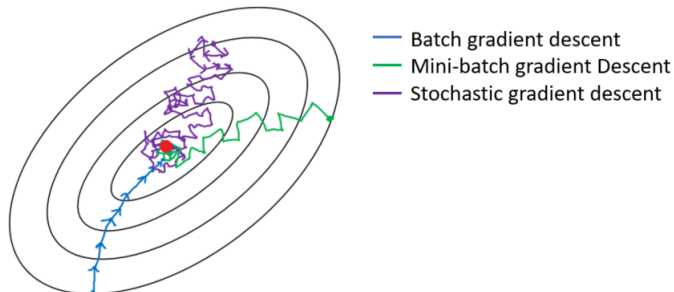
Conceptual Assignment 2

We will look at the numerical impact of ℓ_1 and ℓ_2 norms (used in Lasso and Ridge Regression) on the weights learned in one of the conceptual assignments.

Algorithmic foundations to Machine Learning

Underlying Engine behind ML Training

(Mini-batch) Stochastic Gradient Descent Almost every model and problem-space in ML uses SGD of some kind - Clustering, Regression, Deep Learning, Computer Vision and NLP to name a few. Almost every algorithm in every library - Scikit-learn, Keras, Pytorch, etc uses **mini-batch SGD under the hood**.



So what is Gradient Descent?

Fundamentally

Take a convex/non-convex function, f . GD allows you to find a local optimum to f .

So what is Gradient Descent?

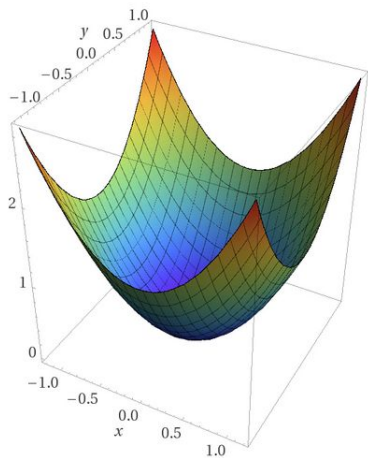
Fundamentally

Take a convex/non-convex function, f . GD allows you to find a local optimum to f .

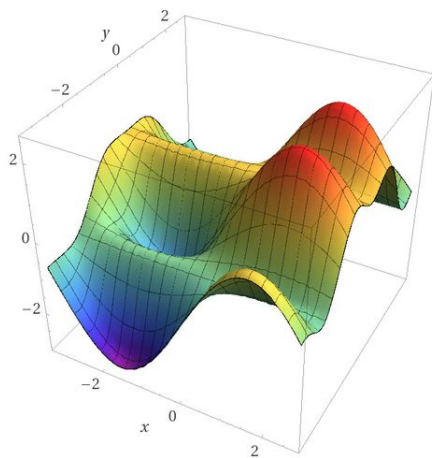
Why is this important?

Consider the Linear Regression problem. \hat{w} is a local optimum to the function $f(w) = \frac{1}{2} \|Xw - y\|_2^2 + \lambda \|w\|_2^2$

Negative Gradient helps you view the direction of descent



Computed by Wolfram|Alpha



Computed by Wolfram|Alpha

Negative Gradients on a Kauai peak!



Gradient Descent

Batch Gradient Descent

Let us say we want to minimize $L(w)$ - Loss Function and find the best \hat{w} that does that.

- 1 **Initialize** $w = w_0$ (maybe randomize)

Gradient Descent

Batch Gradient Descent

Let us say we want to minimize $L(w)$ - Loss Function and find the best \hat{w} that does that.

- 1 **Initialize** $w = w_0$ (maybe randomize)
- 2 **Gradient Descent** $w \leftarrow w - lr * \nabla L(w)$

Gradient Descent

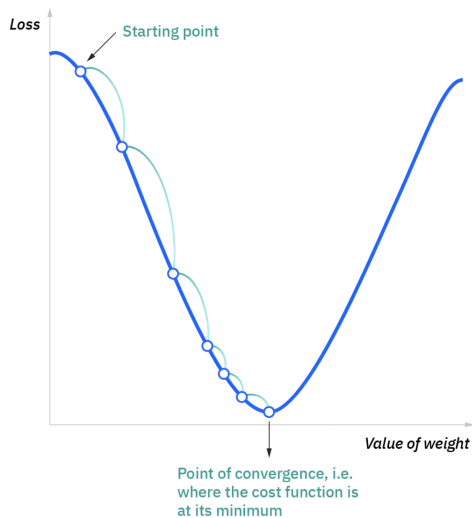
Batch Gradient Descent

Let us say we want to minimize $L(w)$ - Loss Function and find the best \hat{w} that does that.

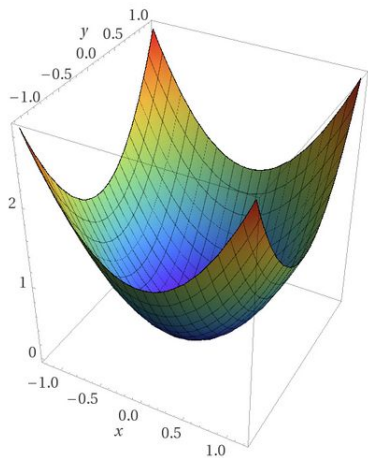
- 1 **Initialize** $w = w_0$ (maybe randomize)
- 2 **Gradient Descent** $w \leftarrow w - lr * \nabla L(w)$
- 3 **Iterate** Repeat step 2 until w converges, i.e.

$$\|w^{k+1} - w^k\| / \|w^k\| \leq 10^{-3}$$

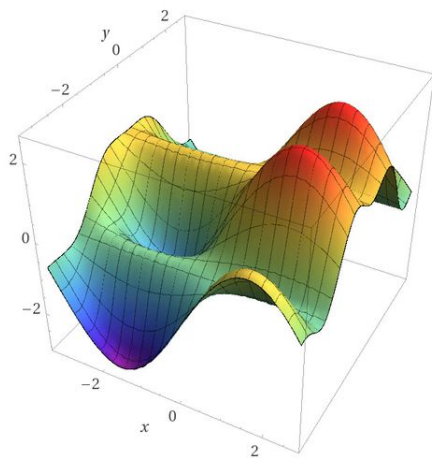
GD in one dimension



Loss function in 2 dimensions



Computed by Wolfram|Alpha



Computed by Wolfram|Alpha

ICE #3

Gradient of Ridge Regularizer (2 mins)

Find the gradient of the regularization function, $R(w) = \lambda \|w\|_2^2$. I.e. obtain the expression for, $\nabla_w R(w)$?

- a) $2\lambda \|w\|_2$
- b) $\lambda \|w\|_2 w$
- c) $2\lambda w$
- d) $2\lambda \|w\|_2 w$

ICE #3

Gradient of Ridge Regularizer (2 mins)

Find the gradient of the regularization function, $R(w) = \lambda \|w\|_2^2$. I.e. obtain the expression for, $\nabla_w R(w)$?

- a) $2\lambda \|w\|_2$
- b) $\lambda \|w\|_2 w$
- c) $2\lambda w$
- d) $2\lambda \|w\|_2 w$

In Assignment 2

We will have a question comparing GD and exact solution for Ridge Regression! Comparison on computation time and accuracy and how both the methods scale?

Gradient Descent Properties

- ① Gradient Descent converges to a local minimum

Gradient Descent Properties

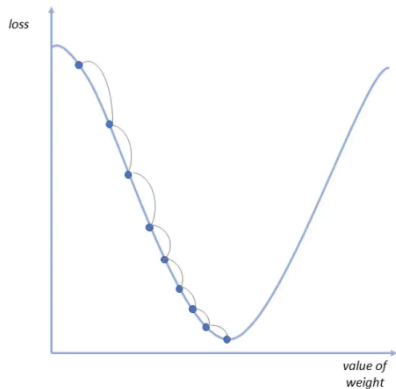
- ① Gradient Descent converges to a local minimum
- ② If L is a convex function, all local minima become a global minima!

Gradient Descent Properties

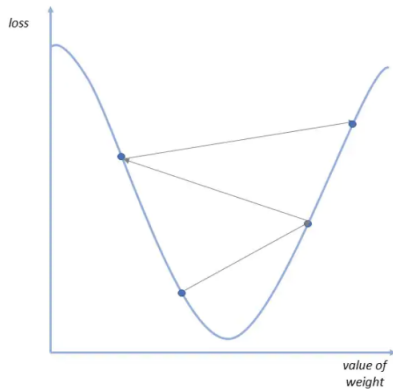
- 1 Gradient Descent converges to a local minimum
- 2 If L is a convex function, all local minima become a global minima!
- 3 Wherever we start, gradient descent usually finds a local minima closest to the start.

Effect of Learning Rate

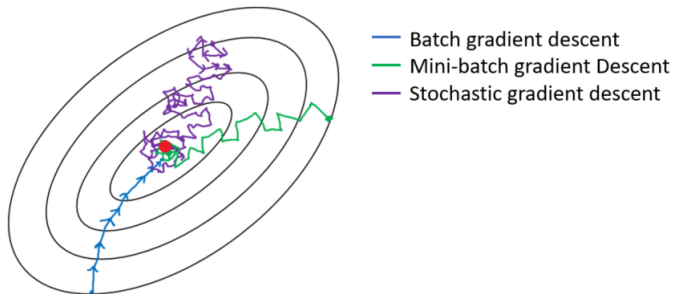
Small Learning Rate



Large Learning Rate



GD behavior in the search space



Gradient descent in practice - SGD!

SGD

Let $L(w) = \sum_{i=1}^N L_i(w)$ where L_i is a function of only the i th data point (x_i, y_i) and parameter w .

- 1 **Initialize** w^0 (randomize)

Gradient descent in practice - SGD!

SGD

Let $L(w) = \sum_{i=1}^N L_i(w)$ where L_i is a function of only the i th data point (x_i, y_i) and parameter w .

- 1 **Initialize** w^0 (randomize) Pick index i at random between 1 and N !

Gradient descent in practice - SGD!

SGD

Let $L(w) = \sum_{i=1}^N L_i(w)$ where L_i is a function of only the i th data point (x_i, y_i) and parameter w .

- 1 **Initialize** w^0 (randomize) Pick index i at random between 1 and N !
- 2 **Gradient Descent** $w^{k+1} \leftarrow w - lr * \nabla L_i(w^k)$

Gradient descent in practice - SGD!

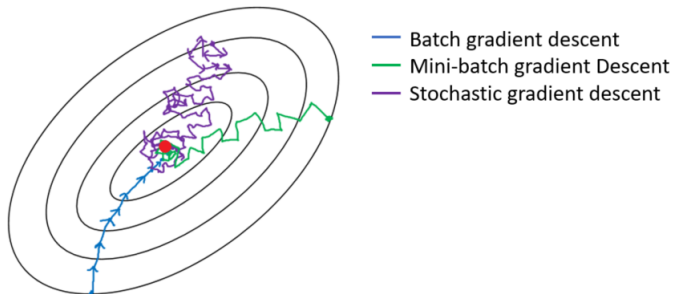
SGD

Let $L(w) = \sum_{i=1}^N L_i(w)$ where L_i is a function of only the i th data point (x_i, y_i) and parameter w .

- 1 **Initialize** w^0 (randomize) Pick index i at random between 1 and N !
- 2 **Gradient Descent** $w^{k+1} \leftarrow w - lr * \nabla L_i(w^k)$
- 3 **Iterate** Repeat step 2 and 3 until w converges, i.e.

$$\|w^{k+1} - w^k\| / \|w^k\| \leq 10^{-3}$$

SGD behavior in search space



SGD in practice - mini-batch SGD!

mini-batch SGD

Let $L(w) = \sum_{i=1}^N L_i(w)$ where L_i is a function of only the i th data point (x_i, y_i) and parameter w . Let B be the number of batches and k be the batch size.

- 1 **Initialize** $w = w_0$ (randomize)

SGD in practice - mini-batch SGD!

mini-batch SGD

Let $L(w) = \sum_{i=1}^N L_i(w)$ where L_i is a function of only the i th data point (x_i, y_i) and parameter w . Let B be the number of batches and k be the batch size.

- 1 **Initialize** $w = w_0$ (randomize) Pick a batch of k data points at random between 1 and N : $i_1, i_2, \dots, i_k!$

SGD in practice - mini-batch SGD!

mini-batch SGD

Let $L(w) = \sum_{i=1}^N L_i(w)$ where L_i is a function of only the i th data point (x_i, y_i) and parameter w . Let B be the number of batches and k be the batch size.

- 1 **Initialize** $w = w_0$ (randomize) Pick a batch of k data points at random between 1 and N : $i_1, i_2, \dots, i_k!$
- 2 **Gradient Descent** $w^{k+1} \leftarrow w^k - lr * \sum_{j=1}^k \nabla_w L_{i_j}(w^k)$

SGD in practice - mini-batch SGD!

mini-batch SGD

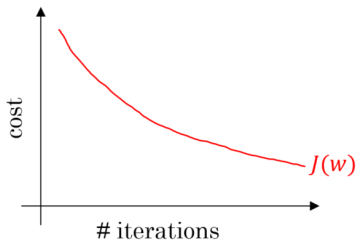
Let $L(w) = \sum_{i=1}^N L_i(w)$ where L_i is a function of only the i th data point (x_i, y_i) and parameter w . Let B be the number of batches and k be the batch size.

- 1 **Initialize** $w = w_0$ (randomize) Pick a batch of k data points at random between 1 and N : $i_1, i_2, \dots, i_k!$
- 2 **Gradient Descent** $w^{k+1} \leftarrow w^k - lr * \sum_{j=1}^k \nabla_w L_{i_j}(w^k)$
- 3 **Iterate** Repeat step 2 and 3 until w converges, i.e.

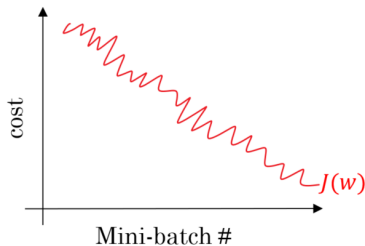
$$\|w^{k+1} - w^k\| / \|w^k\| \leq 10^{-3}$$

GD vs Mini-batch convergence behavior

Batch gradient descent



Mini-batch gradient descent



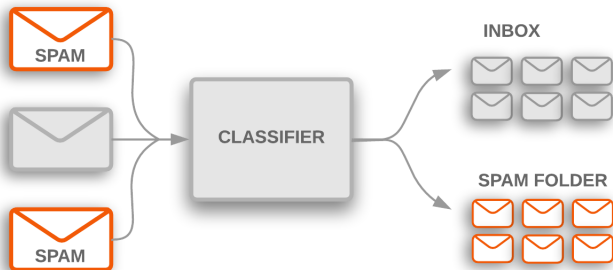
GD vs mini-batch SGD

Factor	GD	Mini-batch SGD
Data	All per iteration	Mini-batch (usually 128 or 256)
Randomness	Deterministic	Stochastic
Error reduction	Monotonic	Stochastic
Computation	High	Low
Memory big data	Intractable	Tractable
Convergence	Low relative error	Few “passes” on data
Local Minima traps	Yes	No

Course Outline

Week	Lecture Material	Assignment
1	Linear Regression	Housing Price Prediction
2	Classification	Spam classification (Kaggle)
3	Classification	Flower/Leaf classification
4	Clustering	MNIST digits clustering
5	Anomaly Detection	Crypto Prediction (Kaggle + P)
6	Data Visualization	Crypto Prediction (Kaggle + P)
7	Deep Learning	Visualizing 1000 images
8	Deep Learning (DL)	ECG Arrhythmia Detection
9	DL in NLP	TwitterSentiment Analysis (Kaggle + P)
10	DLs in Vision	TwitterSentiment Analysis (Kaggle + P)

Classification in Machine Learning



Difference between Classification and Regression

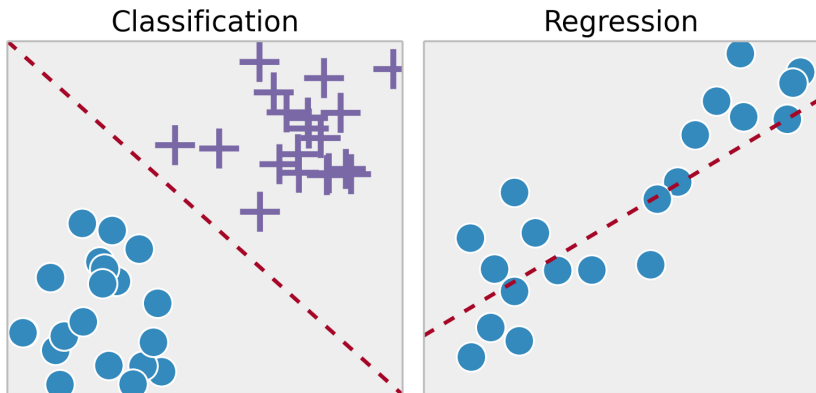
Simple difference

The target type in Regression is **numeric** whereas that in classification is **categorical**

Difference between Classification and Regression

Simple difference

The target type in Regression is **numeric** whereas that in classification is **categorical**



Types of Classification

Binary vs Multi-class classification

With binary categories, its a binary classification problem and with multiple categories, we have a multi-class classification.

Types of Classification

Binary vs Multi-class classification

With binary categories, its a binary classification problem and with multiple categories, we have a multi-class classification.

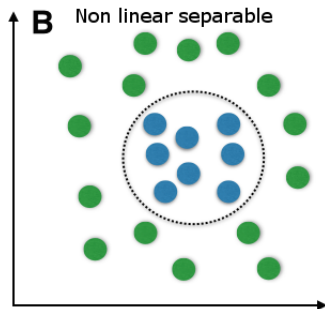
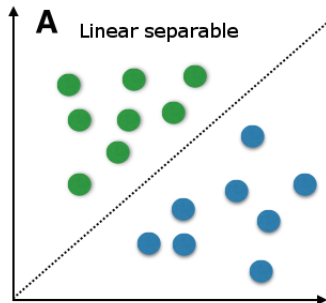
Target is called Label

For binary classification, the convention is to label the target as positive or negative. Example: Positive for spam and negative for not-spam

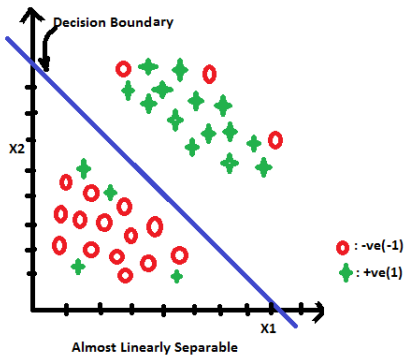
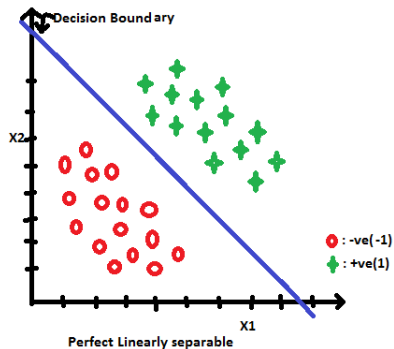
Spam Classification Example

Email excerpt	Type	Label
Could you please respond by tomorrow?	Not-spam	-1
Congratulations!!! You have been selected...	Spam	+1
Looking forward to your presentation...	Not-spam	-1
...

Linear Separability



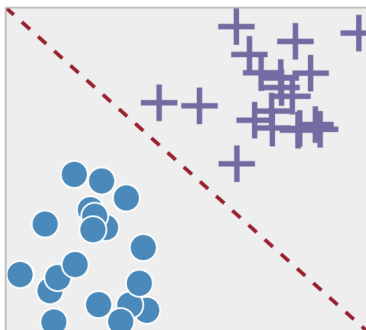
Approximate Linear Separability



ICE #4

Which of the following data sets is the closest to being linearly separable?

Logistic Regression



LR fundamentals

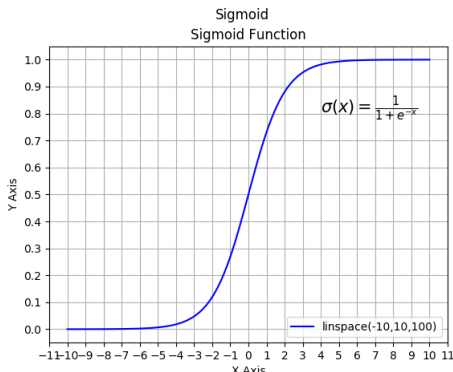
- Linear Model
- Want score $w^T x^i > 0$ for $y_i = +1$ and $w^T x_i < 0$ for $y_i = -1$!
- If linearly separable data, above is feasible. Else, minimize error in separability!!

Logistic Regression

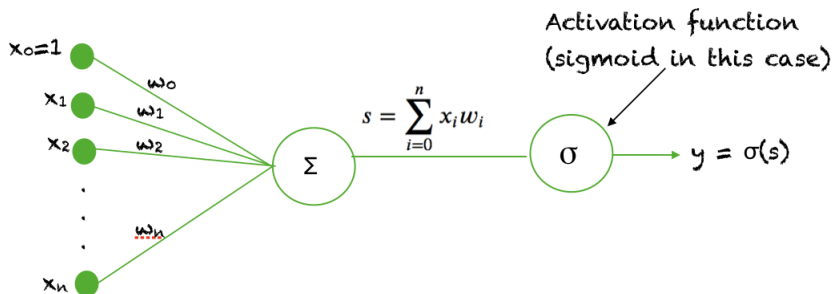
Probability for a class

In LR, the score, $w^T x$ is converted to a probability through the sigmoid function. So we can talk about $P(\hat{y}^i = +1)$ or $P(\hat{y}^i = -1)$

Sigmoid Function



LR represented Graphically



Logistic Regression

LR Prediction

$$\hat{y}_i = \frac{1}{1 + e^{-\hat{w}^T x^i}}$$

LR Loss

Assume that $y_i = 0$ or $y_i = 1$ (i.e. the negative class has a label 0). Then the binary cross-entropy loss applies to LR:

$$\min_w y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)$$

Summary

- Why gradients are important?
- GD vs SGD vs Mini-batch SGD
- Why mini-batch SGD is preferred?
- Regression vs Classification
- Decision Boundary and Linear Separability
- Logistic Regression