

CSCI567 Machine Learning (Spring 2021)

Sirisha Rambhatla

University of Southern California

Jan 29, 2021

1 / 29

Review of Last Lecture

Outline

- 1 Review of Last Lecture
- 2 Linear Classifier and Surrogate Losses
- 3 Perceptron

3 / 29

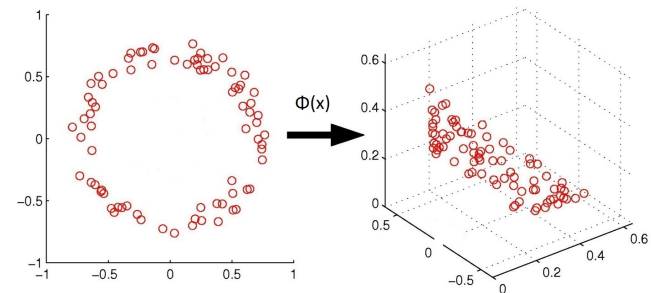
Outline

- 1 Review of Last Lecture
- 2 Linear Classifier and Surrogate Losses
- 3 Perceptron

2 / 29

Review of Last Lecture

Regression with nonlinear basis



Model: $f(x) = w^T \phi(x)$ where $w \in \mathbb{R}^M$

Similar least square solution: $w^* = (\Phi^T \Phi)^{-1} \Phi^T y$

4 / 29

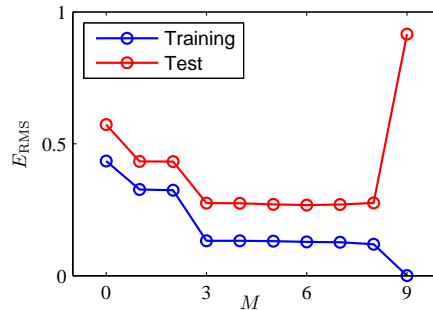
Underfitting and Overfitting

$M \leq 2$ is **underfitting** the data

- large training error
- large test error

$M \geq 9$ is **overfitting** the data

- small training error
- **large test error**



How to prevent overfitting? more data + regularization

$$\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmin}} (\operatorname{RSS}(\mathbf{w}) + \lambda \|\mathbf{w}\|_2^2) = (\Phi^T \Phi + \lambda \mathbf{I})^{-1} \Phi^T \mathbf{y}$$

General idea to derive ML algorithms

Step 1. Pick a set of **models** \mathcal{F}

- e.g. $\mathcal{F} = \{f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} \mid \mathbf{w} \in \mathbb{R}^D\}$
- e.g. $\mathcal{F} = \{f(\mathbf{x}) = \mathbf{w}^T \Phi(\mathbf{x}) \mid \mathbf{w} \in \mathbb{R}^M\}$

Step 2. Define **error/loss** $L(y', y)$

Step 3. Find **empirical risk minimizer (ERM)**:

$$\mathbf{f}^* = \underset{f \in \mathcal{F}}{\operatorname{argmin}} \sum_{n=1}^N L(f(x_n), y_n)$$

or **regularized empirical risk minimizer**:

$$\mathbf{f}^* = \underset{f \in \mathcal{F}}{\operatorname{argmin}} \sum_{n=1}^N L(f(x_n), y_n) + \lambda R(f)$$

ML becomes optimization

Outline

- 1 Review of Last Lecture
- 2 **Linear Classifier and Surrogate Losses**
- 3 Perceptron

Classification

Recall the setup:

- input (feature vector): $\mathbf{x} \in \mathbb{R}^D$
- output (label): $y \in [C] = \{1, 2, \dots, C\}$
- goal: learn a mapping $f: \mathbb{R}^D \rightarrow [C]$

This lecture: **binary classification**

- Number of classes: $C = 2$
- Labels: $\{-1, +1\}$ (cat or dog, fraud or not, price up or down...)

We have discussed **nearest neighbor classifier**:

- require carrying the training set
- more like a heuristic

Deriving classification algorithms

Let's follow the recipe:

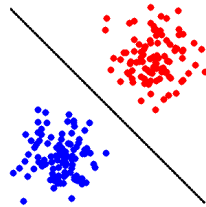
Step 1. Pick a set of models \mathcal{F} .

Again try linear models, but how to predict a label using $w^T x$?

Sign of $w^T x$ predicts the label:

$$\text{sign}(w^T x) = \begin{cases} +1 & \text{if } w^T x > 0 \\ -1 & \text{if } w^T x \leq 0 \end{cases}$$

(Sometimes use sgn for sign too.)



9 / 29

The models

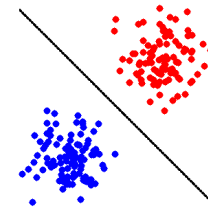
The set of **(separating) hyperplanes**:

$$\mathcal{F} = \{f(x) = \text{sgn}(w^T x) \mid w \in \mathbb{R}^D\}$$

Good choice for **linearly separable** data, i.e., $\exists w$ s.t.

$$\text{sgn}(w^T x_n) = y_n \quad \text{or} \quad y_n w^T x_n > 0$$

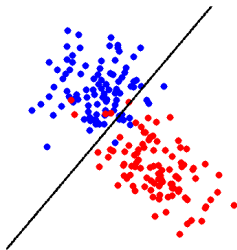
for all $n \in [N]$.



10 / 29

The models

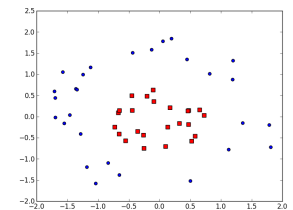
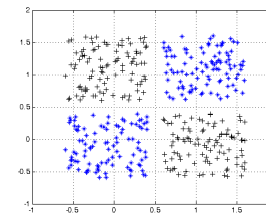
Still makes sense for “almost” linearly separable data



11 / 29

The models

For clearly not linearly separable data,



Again can apply a **nonlinear mapping** Φ :

$$\mathcal{F} = \{f(x) = \text{sgn}(w^T \Phi(x)) \mid w \in \mathbb{R}^M\}$$

More discussions in the next two lectures.

12 / 29

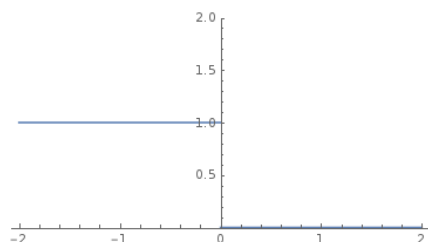
0-1 Loss

Step 2. Define error/loss $L(y', y)$.

Most natural one for classification: **0-1 loss** $L(y', y) = \mathbb{I}[y' \neq y]$

For classification, more convenient to look at the loss **as a function of $yw^T x$** (see ESL 4.5). That is, with

$$\ell_{0-1}(z) = \mathbb{I}[z \leq 0]$$

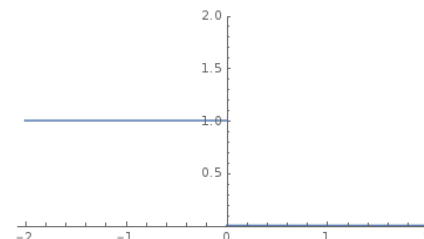


the loss for hyperplane w on example (x, y) is $\ell_{0-1}(yw^T x)$

13 / 29

Minimizing 0-1 loss is hard

However, 0-1 loss is *not convex*.

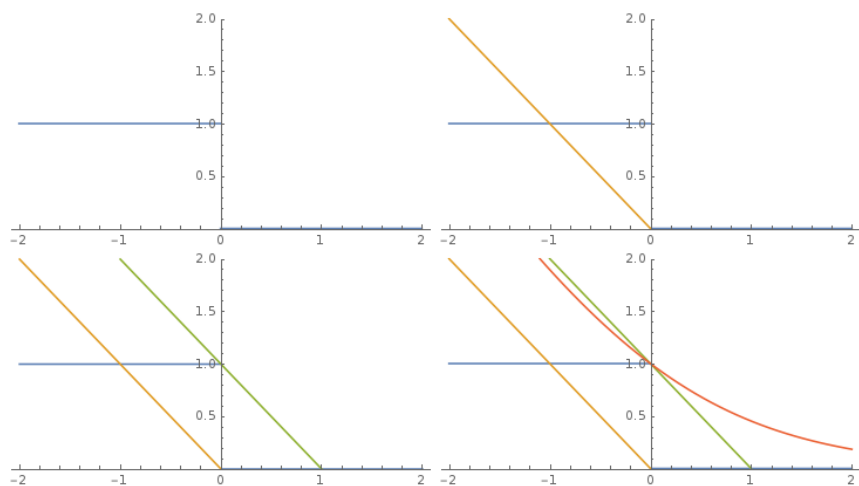


Even worse, minimizing 0-1 loss is *NP-hard in general*.

14 / 29

Surrogate Losses

Solution: find a **convex surrogate loss**



• **perceptron loss** $\ell_{\text{perceptron}}(z) = \max\{0, -z\}$ (used in Perceptron)

• **hinge loss** $\ell_{\text{hinge}}(z) = \max\{0, 1 - z\}$ (used in SVM and many others)

15 / 29

ML becomes convex optimization

Step 3. Find ERM:

$$w^* = \operatorname{argmin}_{w \in \mathbb{R}^D} \sum_{n=1}^N \ell(y_n w^T x_n) = \operatorname{argmin}_{w \in \mathbb{R}^D} \frac{1}{N} \sum_{n=1}^N \ell(y_n w^T x_n)$$

where $\ell(\cdot)$ can be perceptron/hinge/logistic loss

- *no closed-form* in general (unlike linear regression)
- can apply general convex optimization methods

Note: minimizing perceptron loss *does not really make sense* (try $w = 0$), but the algorithm derived from this perspective does.

16 / 29

Outline

- 1 Review of Last Lecture
- 2 Linear Classifier and Surrogate Losses
- 3 **Perceptron**
 - Numerical optimization
 - Applying (S)GD to perceptron loss

17 / 29

The Perceptron Algorithm

In one sentence: **Stochastic Gradient Descent** applied to perceptron loss

i.e. find the minimizer of

$$\begin{aligned} F(\mathbf{w}) &= \frac{1}{N} \sum_{n=1}^N \ell_{\text{perceptron}}(y_n \mathbf{w}^T \mathbf{x}_n) \\ &= \frac{1}{N} \sum_{n=1}^N \max\{0, -y_n \mathbf{w}^T \mathbf{x}_n\} \end{aligned}$$

using SGD

18 / 29

A detour of numerical optimization methods

We describe two simple yet extremely popular methods

- **Gradient Descent (GD)**: simple and fundamental
- **Stochastic Gradient Descent (SGD)**: faster, effective for large-scale problems

Gradient is sometimes referred to as *first-order* information of a function. Therefore, these methods are called *first-order methods*.

19 / 29

Gradient Descent (GD)

Goal: minimize $F(\mathbf{w})$

Algorithm: keep moving in the *negative gradient direction*

Start from some $\mathbf{w}^{(0)}$. For $t = 0, 1, 2, \dots$

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \eta \nabla F(\mathbf{w}^{(t)})$$

where $\eta > 0$ is called step size or learning rate

- in theory η should be set in terms of some parameters of F
- in practice we just try several small values

20 / 29

An example

Example: $F(\mathbf{w}) = 0.5(w_1^2 - w_2)^2 + 0.5(w_1 - 1)^2$. Gradient is

$$\frac{\partial F}{\partial w_1} = 2(w_1^2 - w_2)w_1 + w_1 - 1 \quad \frac{\partial F}{\partial w_2} = -(w_1^2 - w_2)$$

GD:

- Initialize $w_1^{(0)}$ and $w_2^{(0)}$ (to be 0 or *randomly*), $t = 0$
- do

$$w_1^{(t+1)} \leftarrow w_1^{(t)} - \eta \left[2(w_1^{(t)^2} - w_2^{(t)})w_1^{(t)} + w_1^{(t)} - 1 \right]$$

$$w_2^{(t+1)} \leftarrow w_2^{(t)} - \eta \left[-(w_1^{(t)^2} - w_2^{(t)}) \right]$$

$$t \leftarrow t + 1$$

- until $F(\mathbf{w}^{(t)})$ **does not change much**

21 / 29

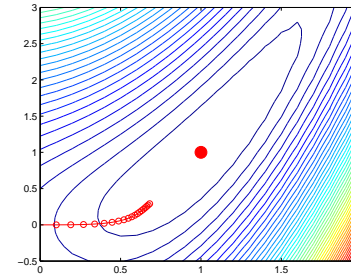
Why GD?

Intuition: by first-order **Taylor approximation**

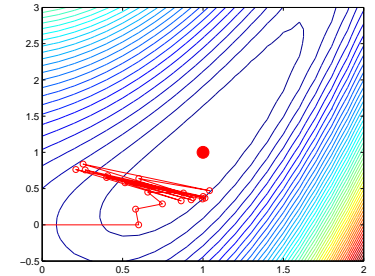
$$F(\mathbf{w}) \approx F(\mathbf{w}^{(t)}) + \nabla F(\mathbf{w}^{(t)})^T (\mathbf{w} - \mathbf{w}^{(t)})$$

GD ensures

$$F(\mathbf{w}^{(t+1)}) \approx F(\mathbf{w}^{(t)}) - \eta \|\nabla F(\mathbf{w}^{(t)})\|_2^2 \leq F(\mathbf{w}^{(t)})$$



reasonable η decreases function value



but large η is unstable

22 / 29

Stochastic Gradient Descent (SGD)

GD: keep moving in the negative gradient direction

SGD: keep moving in some *noisy* negative gradient direction

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \eta \tilde{\nabla} F(\mathbf{w}^{(t)})$$

where $\tilde{\nabla} F(\mathbf{w}^{(t)})$ is a random variable (called **stochastic gradient**) s.t.

$$\mathbb{E} [\tilde{\nabla} F(\mathbf{w}^{(t)})] = \nabla F(\mathbf{w}^{(t)}) \quad (\text{unbiasedness})$$

Key point: it could be *much faster to obtain a stochastic gradient!*

23 / 29

Convergence Guarantees

Many for both GD and SGD on convex objectives.

They tell you at most how many iterations you need to achieve

$$F(\mathbf{w}^{(t)}) - F(\mathbf{w}^*) \leq \epsilon$$

Even for *nonconvex objectives*, many recent works show effectiveness of GD/SGD.

24 / 29

Applying GD to perceptron loss

Objective

$$F(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N \max\{0, -y_n \mathbf{w}^T \mathbf{x}_n\}$$

Gradient (or really *sub-gradient*) is

$$\nabla F(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N -\mathbb{I}[y_n \mathbf{w}^T \mathbf{x}_n \leq 0] y_n \mathbf{x}_n$$

(only misclassified examples contribute to the gradient)

GD update

$$\mathbf{w} \leftarrow \mathbf{w} + \frac{\eta}{N} \sum_{n=1}^N \mathbb{I}[y_n \mathbf{w}^T \mathbf{x}_n \leq 0] y_n \mathbf{x}_n$$

Slow: each update makes one pass of the entire training set!

25 / 29

Applying SGD to perceptron loss

How to construct a stochastic gradient?

One common trick: pick one example $n \in [N]$ uniformly at random, let

$$\tilde{\nabla} F(\mathbf{w}^{(t)}) = -\mathbb{I}[y_n \mathbf{w}^T \mathbf{x}_n \leq 0] y_n \mathbf{x}_n$$

clearly unbiased (convince yourself).

SGD update:

$$\mathbf{w} \leftarrow \mathbf{w} + \eta \mathbb{I}[y_n \mathbf{w}^T \mathbf{x}_n \leq 0] y_n \mathbf{x}_n$$

Fast: each update touches only one data point!

Conveniently, objective of most ML tasks is a *finite sum* (over each training point) and the above trick applies!

Exercise: try SGD to minimize RSS for linear regression.

26 / 29

The Perceptron Algorithm

Perceptron algorithm is SGD with $\eta = 1$ applied to perceptron loss:

Repeat:

- Pick a data point \mathbf{x}_n uniformly at random
- If $\text{sgn}(\mathbf{w}^T \mathbf{x}_n) \neq y_n$

$$\mathbf{w} \leftarrow \mathbf{w} + y_n \mathbf{x}_n$$

Note:

- \mathbf{w} is always a *linear combination* of the training examples

27 / 29

Why does it make sense?

If the current weight \mathbf{w} makes a mistake

$$y_n \mathbf{w}^T \mathbf{x}_n < 0$$

then after the update $\mathbf{w}' = \mathbf{w} + y_n \mathbf{x}_n$ we have

$$y_n \mathbf{w}'^T \mathbf{x}_n = y_n \mathbf{w}^T \mathbf{x}_n + y_n^2 \mathbf{x}_n^T \mathbf{x}_n \geq y_n \mathbf{w}^T \mathbf{x}_n$$

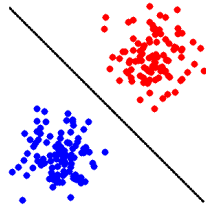
Thus it is more likely to get it right after the update.

28 / 29

Any theory?

(HW 1) If training set is linearly separable

- Perceptron *converges in a finite number of steps*
- training error is 0



There are also guarantees when the data are not linearly separable.