



CIS 419/519 Recitation

Varun Jana, Chang Liu

Oct 13/14 2020

Content

- Perceptron & Activation Functions Overview
- Adagrad
- More about Regularization

Part I: Perceptron & Act. Fns.

and the intuition

Intro to Perceptron

- Our lecture slides have some good content, but let's discuss some roadmaps & high level intuition stuff!
- <https://www.seas.upenn.edu/~cis519/fall2020/assets/lectures/lecture-4/Lecture4-online.pdf> (Page 50 ...)

Intro to Perceptron

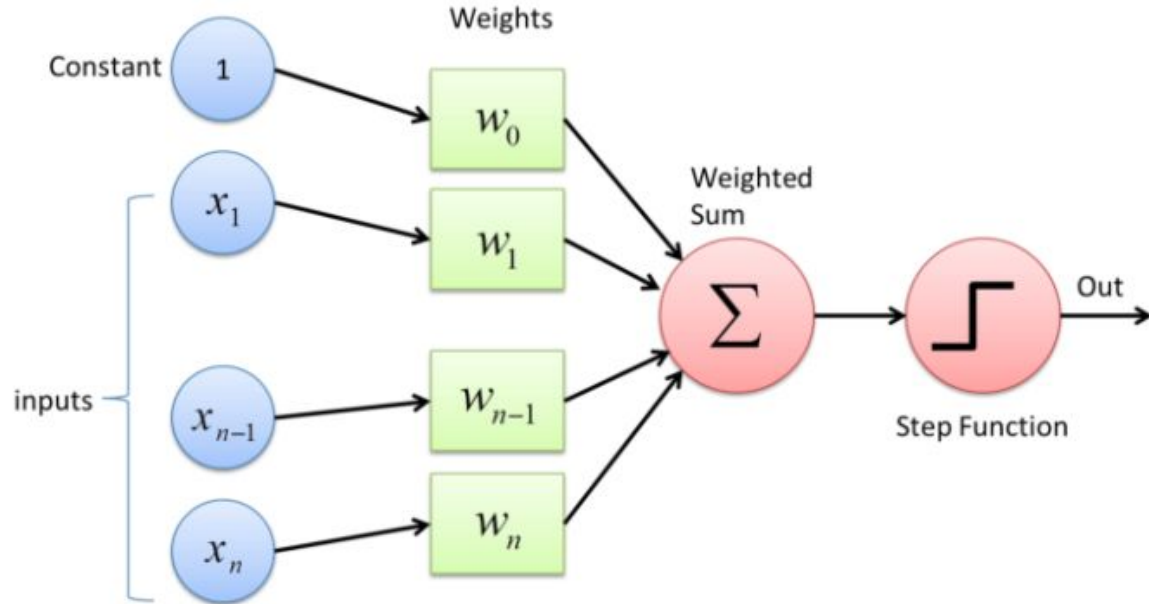


Fig : Perceptron

<https://towardsdatascience.com/what-the-hell-is-perceptron-626217814f53>

Activation Functions - Overview

- Most crucial part of every *neural network* :

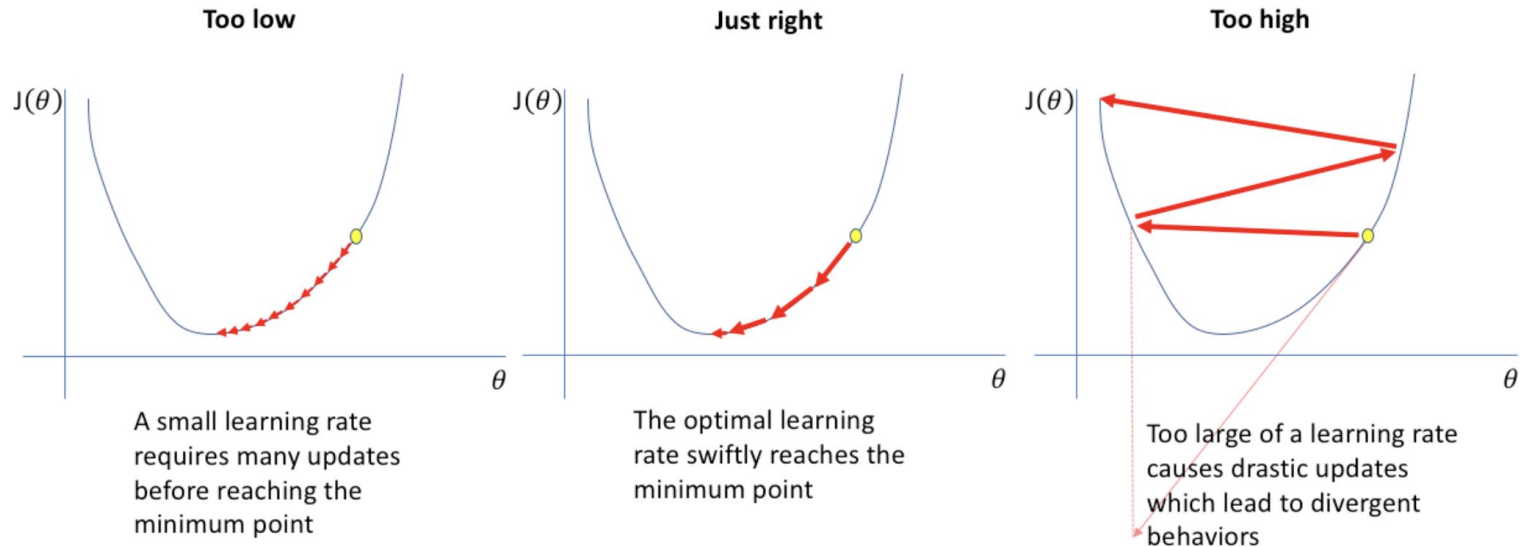
<https://www.analyticssteps.com/blogs/7-types-activation-functions-neural-network>

- Great read! I definitely recommend trying these out in later HWs!

Part 2: Adagrad

What is Adagrad?

- tuning a fix learning rate is tricky
- adaptively scaled the learning rate for each dimension base on historical information

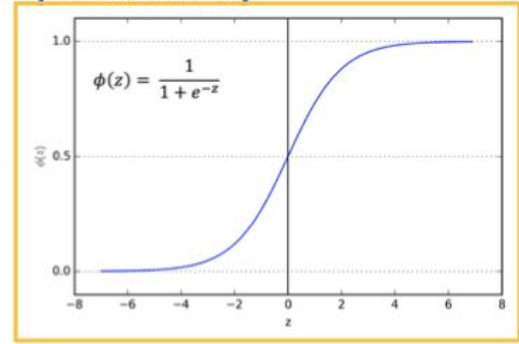


Conditional Probability

- Converting the output of a Perceptron to a conditional probability

$$P(y = +1|\mathbf{x}) = \frac{1}{1 + e^{-Aw^T\mathbf{x}}}$$

- The parameter A can be tuned on a development set



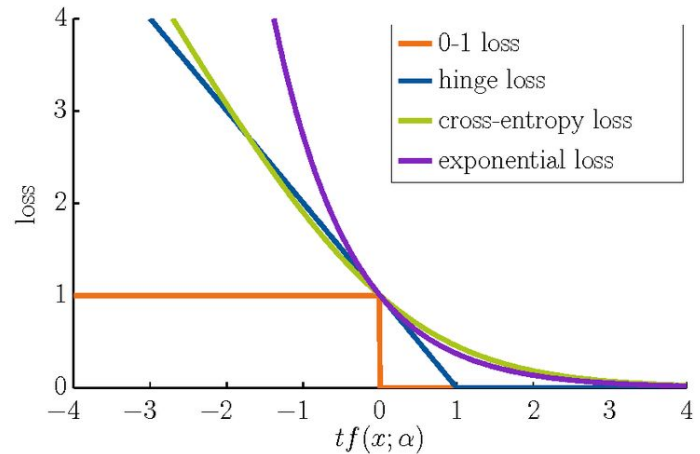
$$\frac{1}{1 + e^{-w^T\mathbf{x}}} > \frac{1}{2}$$

71

Recap on SGD/Gradient Descent

$$\mathbf{w}_{t+1} = \mathbf{w}_t - r_t \mathbf{g}_w \mathcal{L}(\mathbf{x}, y, \mathbf{w}, \theta) = \mathbf{w}_t - r_t \mathbf{g}_t$$

$$\mathcal{L}(\mathbf{x}, y, \mathbf{w}, \theta) = \max \{0, 1 - y(\mathbf{w}^\top \mathbf{x} + \theta)\}$$



Gradient Update

$$\mathbf{w}_{t+1} = \mathbf{w}_t - r_t \mathbf{g}_w Q(\mathbf{z}_t, \mathbf{w}_t) = \mathbf{w}_t - r_t \mathbf{g}_t$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = \begin{cases} \mathbf{0} & \text{if } y(\mathbf{w}^\top \mathbf{x} + \theta) > 1 \\ -y \cdot \mathbf{x} & \text{otherwise} \end{cases}$$
$$\frac{\partial \mathcal{L}}{\partial \theta} = \begin{cases} 0 & \text{if } y(\mathbf{w}^\top \mathbf{x} + \theta) > 1 \\ -y & \text{otherwise} \end{cases}$$

Now! Learning rate!

$$\mathbf{w}_{t+1} = \mathbf{w}_t - r_t \mathbf{g}_w Q(\mathbf{z}_t, \mathbf{w}_t) = \mathbf{w}_t - \boxed{r_t} \mathbf{g}_t$$

$$G_j^t = \sum_{k=1}^t \left(\frac{\partial \mathcal{L}}{\partial w_j^k} \right)^2$$

$$H^t = \sum_{k=1}^t \left(\frac{\partial \mathcal{L}}{\partial \theta^k} \right)^2$$

$$\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t + \eta \cdot \frac{y \cdot \mathbf{x}}{\sqrt{\mathbf{G}^t}}$$

$$\theta^{t+1} \leftarrow \theta^t + \eta \frac{y}{\sqrt{H^t}}$$

Why Adagrad?

Advantages of Using AdaGrad

- it eliminates the need to manually tune the learning rate
- convergence is faster and more reliable – than simple SGD when the scaling of the weights is unequal
- It is not very sensitive to the size of the master step

Overview of Gradient Descent Based Optimization Algorithm:

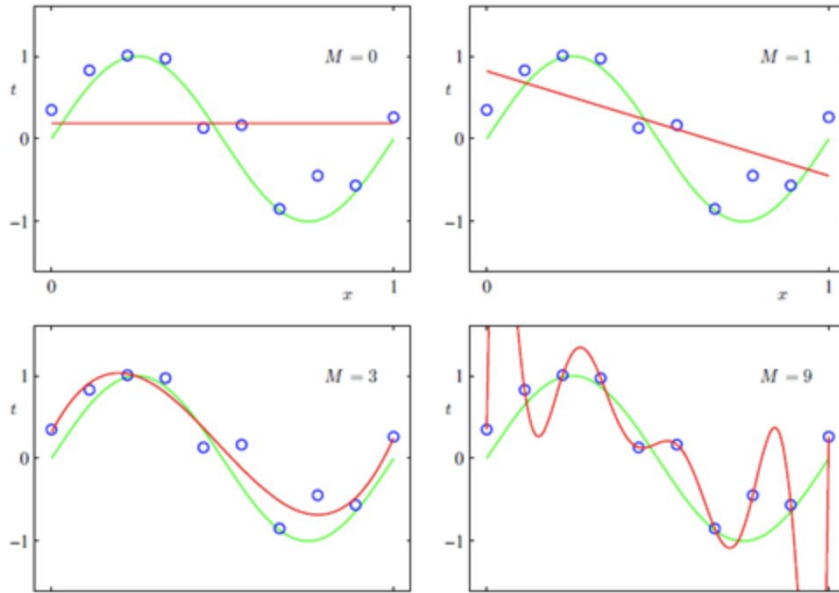
<https://ruder.io/optimizing-gradient-descent/>

Tips for hw

1. separately update w and θ
2. keep track of G for w_i and H for θ
3. use $\sqrt{\text{sum of } G / H}$ to update learning rate

Part 3: Regularization

Why Regularization?



- avoid overfitting
- more robust to noise

Regularization Via Averaged Perceptron

- **Variables:**
 - m : number of examples
 - k : number of mistakes
 - c_i : consistency count for hypothesis v_i
 - T : number of epochs
- **Input:** a labeled training $\{\{\mathbf{x}_1, y_1\}, \{\mathbf{x}_2, y_2\}, \dots, \{\mathbf{x}_m, y_m\}\}$
- **Output:** a list of weighted perceptrons $\{\{\mathbf{v}_1, c_1\}, \dots, \{\mathbf{v}_k, c_k\}\}$
- **Initialize:** $k = 0, \mathbf{v}_1 = \mathbf{0}, c_1 = 0$
- Repeat T times:
 - For $i = 1, \dots, m$;
 - Compute prediction $y' = \text{sgn}(\mathbf{v}_k^T \cdot \mathbf{x}_i)$
 - If $y' \neq y$, then $c_k = c_k + 1$
else: $\mathbf{v}_{k+1} = \mathbf{v}_k + y_i \mathbf{x}; c_{k+1} = 1; k = k + 1$
- **Prediction:**
 - **Given:** a list of weighted perceptrons $\{\{\mathbf{v}_1, c_1\}, \dots, \{\mathbf{v}_k, c_k\}\}$; a new example \mathbf{x}
 - **Predict** the label (\mathbf{x}) as follows: $y(\mathbf{x}) = \text{sgn}[\sum_1^k c_i (\mathbf{v}_i^T \mathbf{x})]$

- This can be done on top of any online mistake driven algorithm.
- In HW 2 you will run it over three different algorithms.
- The implementation requires thinking.

Averaged version of Perceptron/Winnow is as good as any other linear learning algorithm, if not better.

Perceptron

Input set of examples and their labels

$$Z = ((\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)) \in \mathbf{R}^n \times \{-1, 1\}^m, \eta, \theta_{Init}$$

- Initialize $\mathbf{w} \leftarrow 0$ and $\theta \leftarrow \theta_{Init}$
- For every training epoch:
- for every $\mathbf{x}_j \in \mathbf{X}$:
 - $\hat{y} \leftarrow \text{sign}(\langle \mathbf{w}, \mathbf{x}_j \rangle - \theta)$
 - If $(\hat{y} \neq y_j)$
 - $\mathbf{w} \leftarrow \mathbf{w} + \eta y_j \mathbf{x}_j$
 - $\theta \leftarrow \theta + \eta y_j$

Just to make sure we understand that we learn both \mathbf{w} and θ

The Key Ideas

- want to somehow regularize or shrink coefficient (w) towards zero
- In other words, this technique discourages learning a more complex or flexible model, so as to avoid the risk of overfitting

Mathematically

$$J(\mathbf{w}) = \frac{1}{m} \sum_{i=1}^m Q(\mathbf{z}_i, \mathbf{w}_i) + \lambda R_i(\mathbf{w}_i)$$

Loss:

LMS case: $Q((\mathbf{x}, y), \mathbf{w}) = (y - \mathbf{w}^T \mathbf{x})^2$

- $R(\mathbf{w}) = \|\mathbf{w}\|_2^2$ gives the optimization problem called Ridge Regression.
- $R(\mathbf{w}) = \|\mathbf{w}\|_1$ gives a problem called the LASSO problem

Hinge Loss case: $Q((\mathbf{x}, y), \mathbf{w}) = \max(0, 1 - y \mathbf{w}^T \mathbf{x})$

- $R(\mathbf{w}) = \|\mathbf{w}\|_2^2$ gives the problem called Support Vector Machines

Logistics Loss case: $Q((\mathbf{x}, y), \mathbf{w}) = \log(1 + \exp\{-y \mathbf{w}^T \mathbf{x}\})$

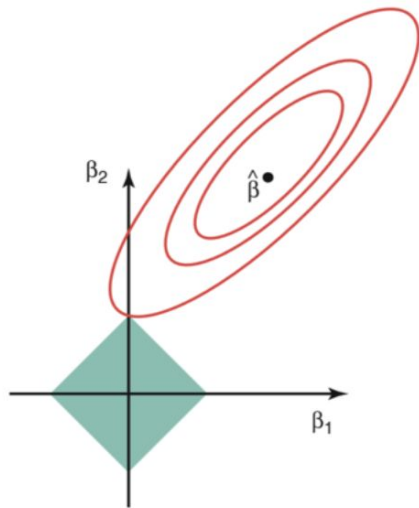
- $R(\mathbf{w}) = \|\mathbf{w}\|_2^2$ gives the problem called Logistics Regression

Regularization:

L2: $\lambda \sum_{i=1}^m |w_i|^2$

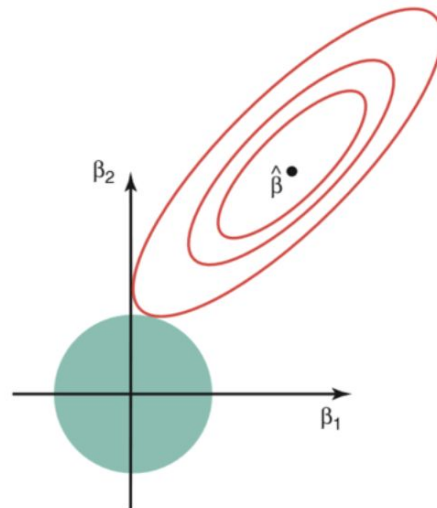
L1: $\lambda \sum_{i=1}^m |w_i|$

Difference between L1 & L2 Regularization



L1:

$$\lambda \sum_{i=1}^m |w_i|$$



L2:

$$\lambda \sum_{i=1}^m |w_i|^2$$

Choose which?

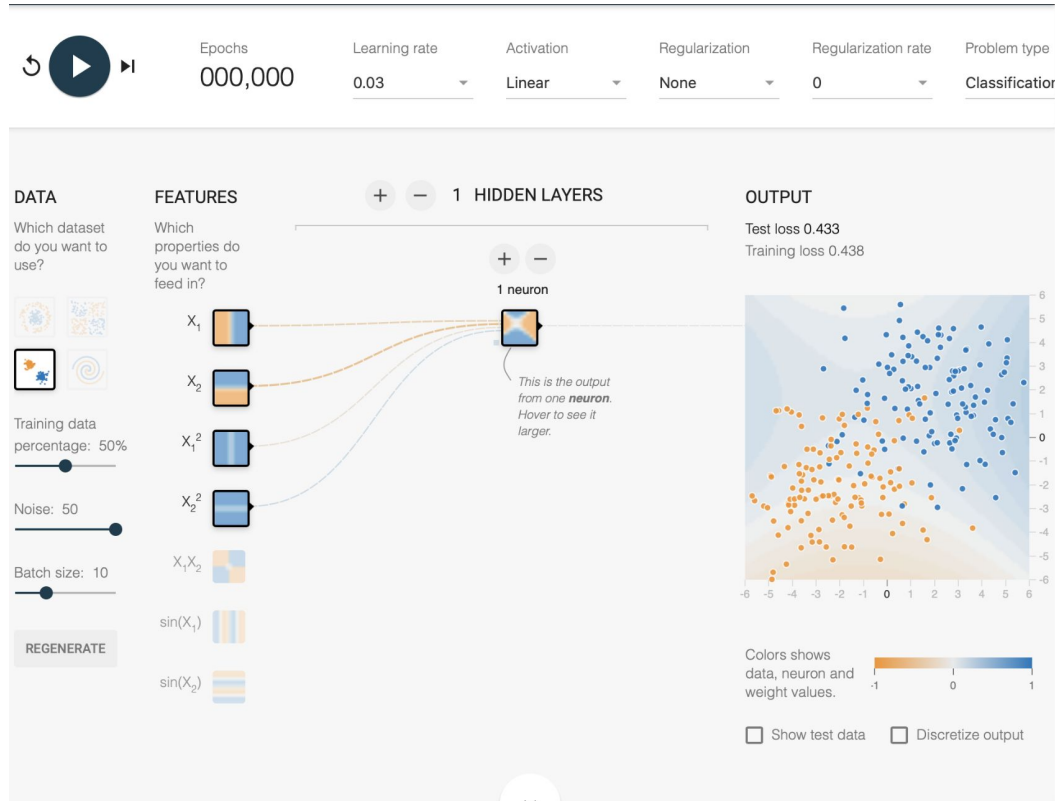
L1:

- penalizes sum of absolute value of weights.
- has a sparse solution
- multiple solutions
- has built in feature selection
- robust to outliers
- generates model that are simple and interpretable but cannot learn complex patterns

L2:

- penalizes sum of square weights.
- has a non sparse solution
- has one solution
- has no feature selection
- is not robust to outliers
- better prediction when output variable is a function of all input features
- is able to learn complex data patterns

Playground



<https://developers.google.com/machine-learning/crash-course/regularization-for-simplicity/playground-exercise-examining-l2-regularization>