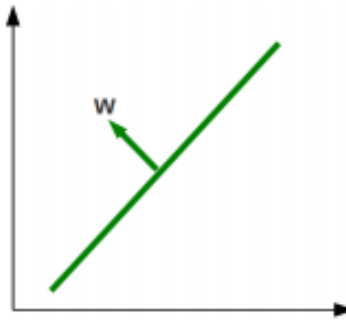# Hyperplane based Classifiers (1): The Perceptron Algorithm

CS771: Introduction to Machine Learning

Piyush Rai

# Hyperplane

- Separates a $D$-dimensional space into two **half-spaces** (positive and negative)
- Defined by a normal vector $\boldsymbol{w} \in \mathbb{R}^D$ (pointing towards positive half-space)
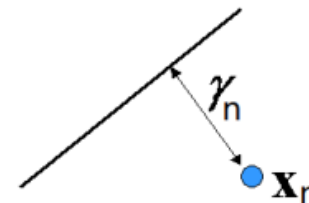
$b > 0$ means moving $\boldsymbol{w}^\top \boldsymbol{x} = 0$ along the direction of $\boldsymbol{w}$; $b < 0$ means in opp. dir.

$$\boldsymbol{w}^\top \boldsymbol{x} + b = 0$$

- Equation of the hyperplane: $\boldsymbol{w}^\top \boldsymbol{x} = 0$
- Assumption: The hyperplane passes through origin. If not, add a bias term $b$
- Distance of a point $\boldsymbol{x}_n$ from a hyperplane $\boldsymbol{w}^\top \boldsymbol{x} + b = 0$
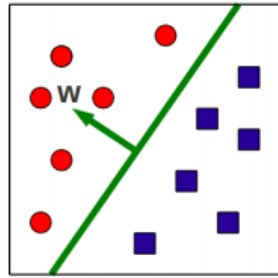
Can be positive or negative

$$\gamma_n = \frac{\boldsymbol{w}^T \boldsymbol{x}_n + b}{||\boldsymbol{w}||}$$

# Hyperplane based (binary) classification

- Basic idea: Learn to separate two classes by a hyperplane $\boldsymbol{w}^\mathsf{T}\boldsymbol{x} + b = 0$
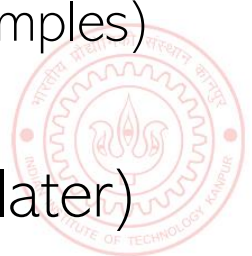
Prediction Rule

$$y_* = \text{sign}(\boldsymbol{w}^\mathsf{T}\boldsymbol{x}_* + b)$$

For multi-class classification with hyperplanes, there will be multiple hyperplanes (e.g., one for each pair of classes); more on this later

- The hyperplane may be "implied" by the model, or learned directly
  - Implied: Prototype-based classification, nearest neighbors, generative classification, etc
  - Directly learned: Logistic regression, Perceptron, Support Vector Machine (SVM), etc
- The "direct" approach defines a model with params $\boldsymbol{w}$ (and optionally a bias param $b$)
  - The parameters are learned by optimizing a classification loss function (will soon see examples)
  - These are also discriminative approaches – $\boldsymbol{x}$ is not modeled but treated as fixed (given)
- The hyperplane need not be linear (e.g., can be made nonlinear using kernels; later)

# Interlude: Loss Functions for Classification

- In regression (assuming linear model $\hat{y} = \boldsymbol{w}^\top \boldsymbol{x}$), some common loss fn
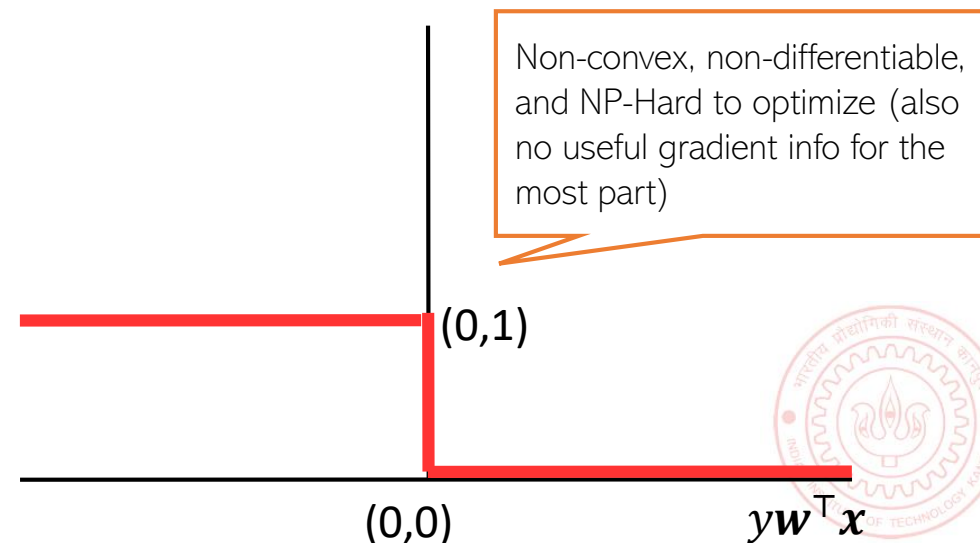
$$\ell(y, \hat{y}) = (y - \hat{y})^2 \qquad\qquad \ell(y, \hat{y}) = |y - \hat{y}|$$

- These measure the difference between the true output and model's prediction

- What about loss functions for <u>classification</u> where $\hat{y} = \text{sign}(\boldsymbol{w}^\top \boldsymbol{x})$ ?

- Perhaps the most natural classification loss function would be a **"0-1 Loss"**

  - Loss = 1 if $\hat{y} \neq y$ and Loss = 0 if $\hat{y} = y$.

  - Assuming labels as +1/-1, it means

$$\ell(y, \hat{y}) = \begin{cases} 1 & \text{if } y\boldsymbol{w}^\top \boldsymbol{x} < 0 \\ 0 & \text{if } y\boldsymbol{w}^\top \boldsymbol{x} \geq 0 \end{cases}$$

0-1 Loss

Non-convex, non-differentiable, and NP-Hard to optimize (also no useful gradient info for the most part)

(0,1)

(0,0)

$y\boldsymbol{w}^\top \boldsymbol{x}$

Same as $\mathbb{I}[y\boldsymbol{w}^\top \boldsymbol{x} < 0]$ or $\mathbb{I}[\text{sign}(\boldsymbol{w}^\top \boldsymbol{x}) \neq y]$
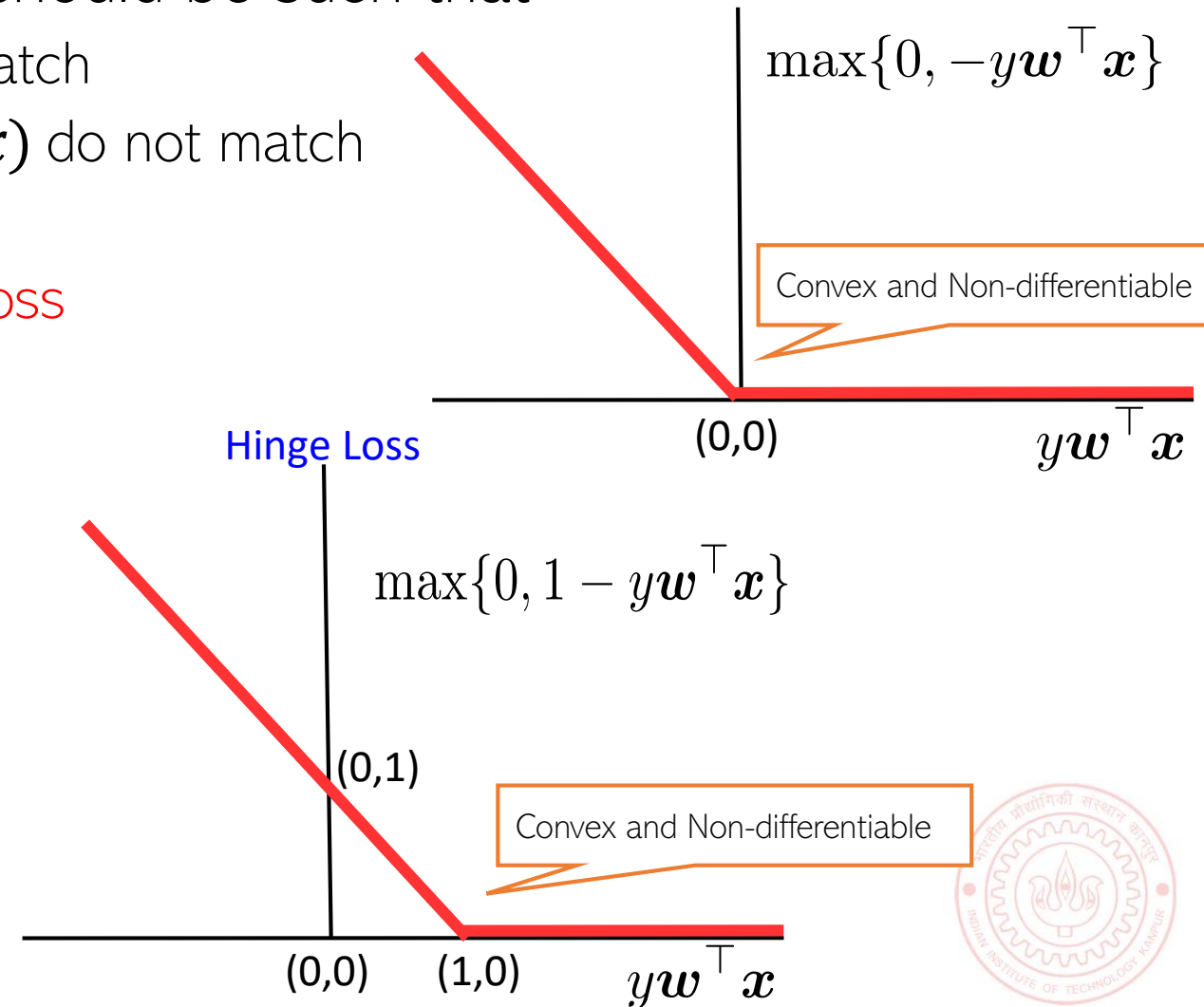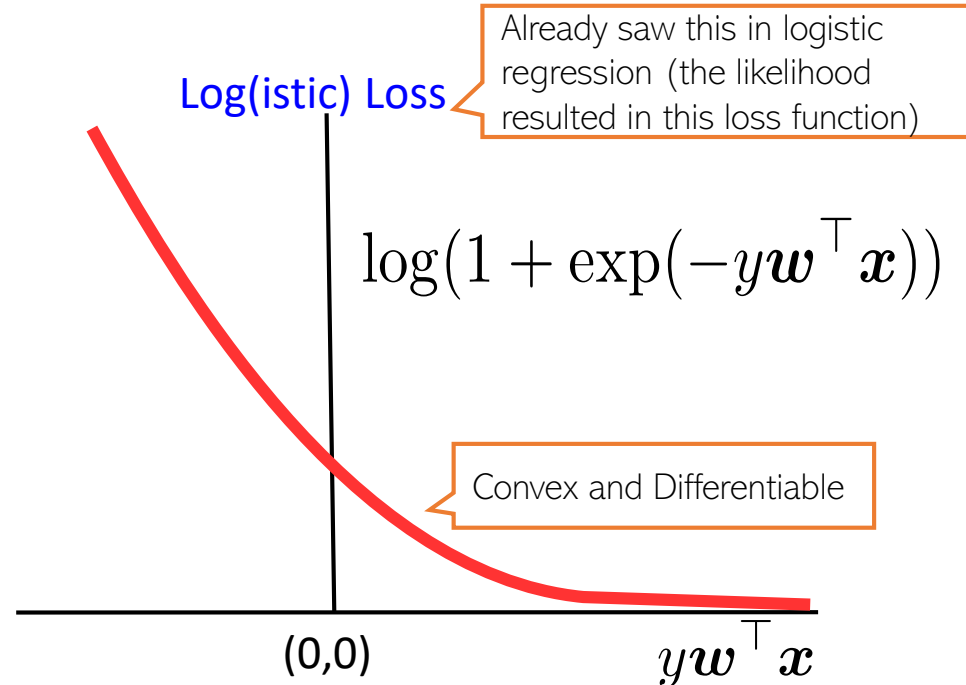
# Interlude: Loss Functions for Classification

- An ideal loss function for classification should be such that
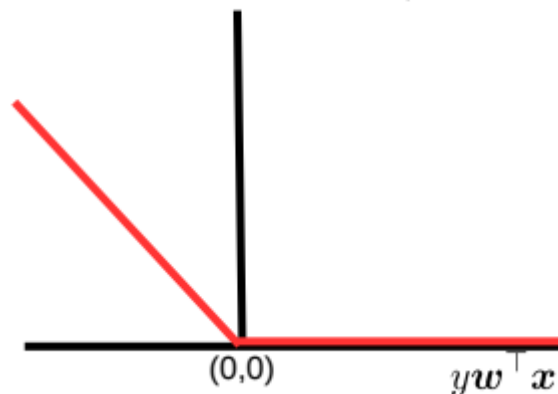  - Loss is small/zero if $y$ and $\text{sign}(w^\top x)$ match
  - Loss is large/non-zero if $y$ and $\text{sign}(w^\top x)$ do not match
  - Large positive $yw^\top x \Rightarrow$ small/zero loss
  - Large negative $yw^\top x \Rightarrow$ large/non-zero loss

**"Perceptron" Loss**

$$\max\{0, -yw^\top x\}$$

Convex and Non-differentiable

(0,0)    $yw^\top x$

Already saw this in logistic regression (the likelihood resulted in this loss function)

**Log(istic) Loss**

$$\log(1 + \exp(-yw^\top x))$$

Convex and Differentiable

(0,0)    $yw^\top x$

**Hinge Loss**

$$\max\{0, 1 - yw^\top x\}$$

(0,1)

Convex and Non-differentiable

(0,0)  (1,0)   $yw^\top x$

# Learning by Optimizing Perceptron Loss

- Let's ignore the bias term $b$ for now. So the hyperplane is simply $\mathbf{w}^\top \mathbf{x} = 0$
- The Perceptron loss function: $L(w) = \sum_{n=1}^{N} \max\{0, -y_n \mathbf{w}^\top \mathbf{x}_n\}$. Let's do SGD

"Perceptron" Loss: $\max\{0, -y\mathbf{w}^\top \mathbf{x}\}$

(0,0)  $y\mathbf{w}^\top \mathbf{x}$

Subgradients w.r.t. $\mathbf{w}$

One randomly chosen example in each iteration

$$\mathbf{g}_n = \begin{cases} 0, & \text{for } y_n \mathbf{w}^\top \mathbf{x}_n > 0 \\ -y_n \mathbf{x}_n & \text{for } y_n \mathbf{w}^\top \mathbf{x}_n < 0 \\ k y_n \mathbf{x}_n & \text{for } y_n \mathbf{w}^\top \mathbf{x}_n = 0 \quad (\text{where } k \in [-1, 0]) \end{cases}$$

- If we use $k = 0$ then $\mathbf{g}_n = 0$ for $y_n \mathbf{w}^\top \mathbf{x}_n \geq 0$, and $\mathbf{g}_n = -y_n \mathbf{x}_n$ for $y_n \mathbf{w}^\top \mathbf{x}_n < 0$
- Non-zero gradients only when the model makes a mistake on current example $(\mathbf{x}_n, y_n)$
- Thus SGD will update $\mathbf{w}$ only when there is a mistake (mistake-driven learning)

# The Perceptron Algorithm

- Stochastic Sub-grad desc on Perceptron loss is also known as the Perceptron algorithm

## Stochastic SubGD

Note: An example may get chosen several times during the entire run

① Initialize $\boldsymbol{w} = \boldsymbol{w}^{(0)}, t = 0$, set $\eta_t = 1, \forall t$

② Pick some $(\boldsymbol{x}_n, y_n)$ randomly.

Mistake condition

③ If current $\boldsymbol{w}$ makes a mistake on $(\boldsymbol{x}_n, y_n)$, i.e., $y_n {\boldsymbol{w}^{(t)}}^\top \boldsymbol{x}_n < 0$

$$\boldsymbol{w}^{(t+1)} = \boldsymbol{w}^{(t)} + y_n \boldsymbol{x}_n$$
$$t = t + 1$$

④ If not converged, go to step 2.

Updates are "corrective": If $y\_n = +1$ and $\boldsymbol{w}^\top \boldsymbol{x}_n < 0$, after the update $\boldsymbol{w}^\top \boldsymbol{x}_n$ will be less negative. Likewise, if $y_n = -1$ and $\boldsymbol{w}^\top \boldsymbol{x}_n > 0$, after the update $\boldsymbol{w}^\top \boldsymbol{x}_n$ will be less positive

If training data is linearly separable, the Perceptron algo will converge in a finite number of iterations (Block & Novikoff theorem)
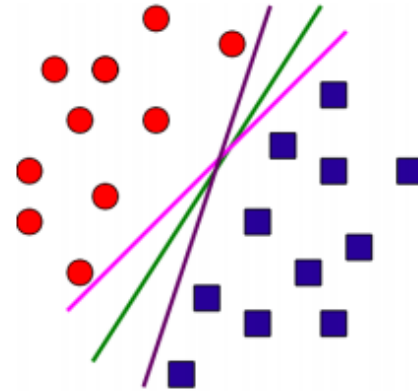
- An example of an online learning algorithm (processes one training ex. at a time)
- Assuming $\boldsymbol{w}^{(0)} = \boldsymbol{0}$, easy to see that the final $\boldsymbol{w}$ has the form $\boldsymbol{w} = \sum_{n=1}^N \alpha_n y_n \boldsymbol{x}_n$

  - $\alpha_n$ is total number of mistakes made by the algorithm on example $(\boldsymbol{x}_n, y_n)$

    Meaning of $\alpha_n$ may be different

  - As we'll see, many other models also have weights $\boldsymbol{w}$ in the form $\boldsymbol{w} = \sum_{n=1}^N \alpha_n y_n \boldsymbol{x}_n$

# Perceptron and (lack of) Margins

- Perceptron would learn a hyperplane (of many possible) that separates the classes



> Basically, it will learn the hyperplane which corresponds to the $\boldsymbol{w}$ that minimizes the Perceptron loss

> Kind of an "unsafe" situation to have – ideally would like it to be reasonably away from closest training examples from either class

- Doesn't guarantee any "margin" around the hyperplane
  - The hyperplane can get arbitrarily close to some training example(s) on either side
  - This may not be good for generalization performance

> $\gamma > 0$ is some pre-specified margin

- Can artificially introduce margin by changing the mistake condition to $y_n \boldsymbol{w}^\top \boldsymbol{x}_n < \gamma$

- Support Vector Machine (SVM) does it directly by learning the max. margin hyperplane

# Coming up next

- Support Vector Machines