

Linear Classifiers, Logistic Regression

CSC 461: Machine Learning

Fall 2022

Prof. Marco Alvarez
University of Rhode Island

Linear classifiers

Linear classifiers

- ▶ Discriminative
 - ✓ Perceptron
 - ✓ Logistic regression
 - ✓ Support vector machines
- ▶ Generative
 - ✓ Linear discriminant analysis
 - ✓ Naive bayes

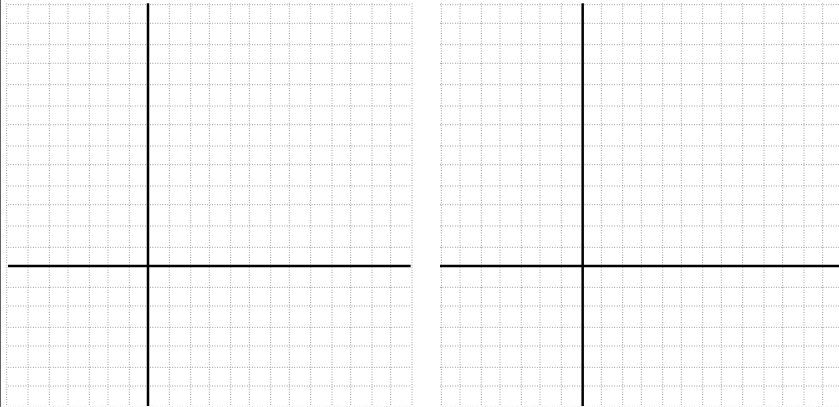
Binary classification

$$\mathcal{D} = \{(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(n)}, y^{(n)})\}$$

$$\mathbf{x}^{(i)} \in \mathbb{R}^d \quad y^{(i)} \in \{-1, +1\}$$

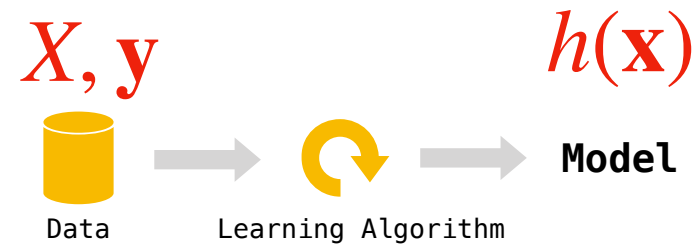
x1	x2	y
0.5	0.1	+1
0.3	0.9	-1
0.3	0.875	-1
0.45	0.15	+1
...

Plots (regression x classification)



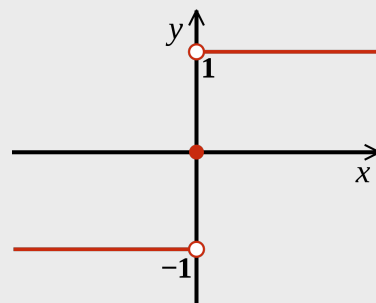
Binary classification goal

- Learn a **decision boundary** such that two classes can be separated



The *sign* function

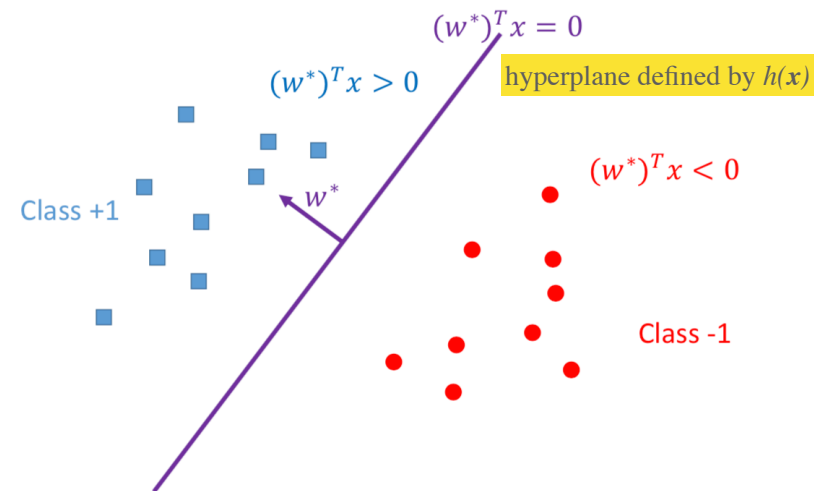
$$\text{sign}(x) = \begin{cases} -1 & \text{if } x < 0, \\ 0 & \text{if } x = 0, \\ +1 & \text{if } x > 0 \end{cases}$$



$$h(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x})$$

Image credit: Wikipedia

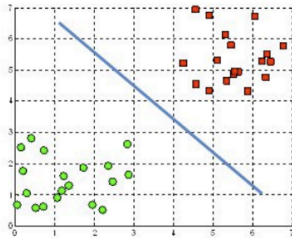
Decision boundary



credit: yingyu liang, cos 495, princeton

Decision boundary

A hyperplane in \mathbb{R}^2 is a line



$$0 = b + w_1x_1 + w_2x_2$$

$$x_2 = -\frac{b}{w_2} - \frac{w_1}{w_2}x_1$$

Image credit: <https://mc.ai/why-activation-function-is-used-in-neural-network/>

Absorbing the bias

$$h(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x} + b)$$

$$= \text{sign}\left(\sum_{i=1}^d w_i x_i + b\right)$$

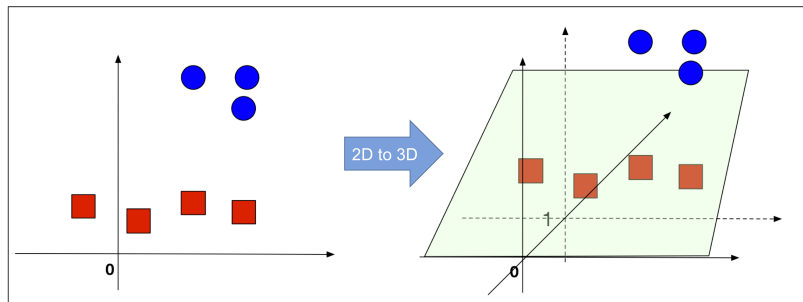
$$x_0 = 1, \quad w_0 = b$$

$$h(\mathbf{x}) = \text{sign}\left(\sum_{i=0}^d w_i x_i\right)$$

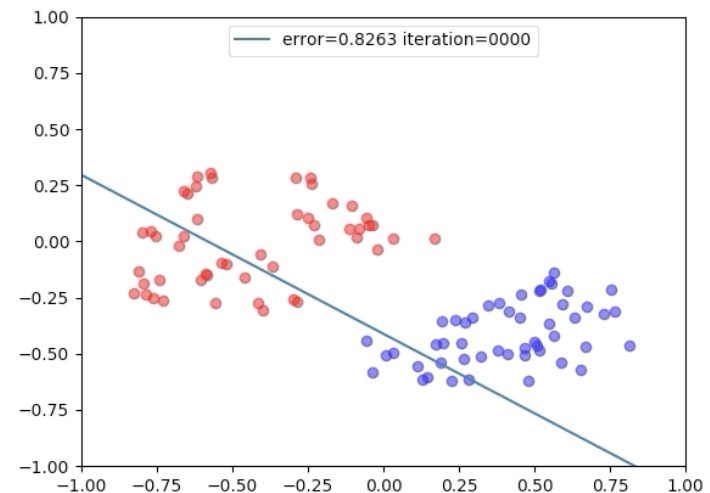
$$= \text{sign}(\mathbf{w}^T \mathbf{x})$$

x0	x1	x2	Y
1	0.5	0.1	+1
1	0.3	0.9	-1
1	0.3	0.875	-1
1	0.25	0.561	-1
1	0.45	0.15	+1
...

Absorbing the bias



Learning



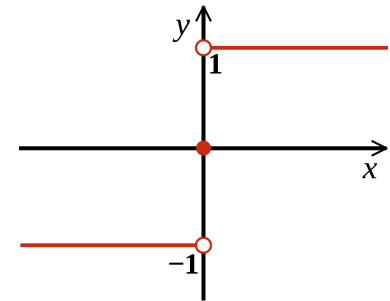
Example

- Provide a solution (weight vector)

x_0	x_1	x_2	y
1	0	0	-1
1	0	1	-1
1	1	0	-1
1	1	1	+1

The *sign* function (again)

$$\text{sign}(x) = \begin{cases} -1 & \text{if } x < 0, \\ 0 & \text{if } x = 0, \\ +1 & \text{if } x > 0 \end{cases}$$

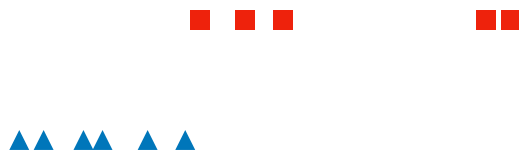


Note that the gradient is zero almost everywhere and the gradient is undefined at $x = 0$.

Image credit: Wikipedia

Can we use the squared loss?

- Treat target labels (binary) as continuous
 - ✓ final prediction decided by checking $h(\mathbf{x}) > 0$



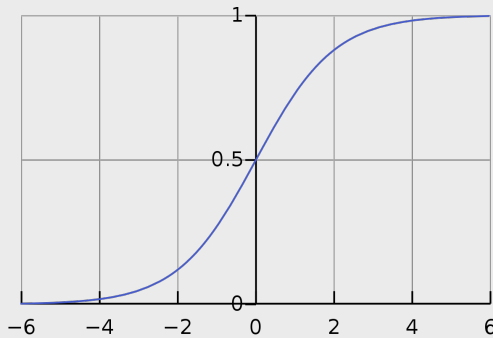
Predicted values can fall outside $[-1, 1]$ range

Square loss penalizes correct predictions with large losses

Logistic regression

Logistic function

$$\sigma(x) = \frac{1}{1 + e^{-x}} = \frac{e^x}{e^x + 1}$$



mapping \mathbb{R} to $[0,1]$

continuous and differentiable

<https://alliance.seas.upenn.edu/~cis520/wiki/index.php?n=Lectures.Logistic>

Logistic regression

Binary classifier

✓ uses a **logistic function** (type of sigmoid function, S-shaped)

✓ models **probability** of output in terms of input

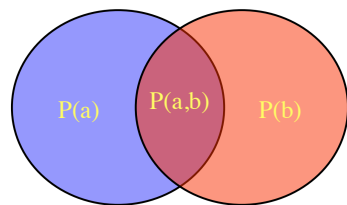
It is considered a **linear classifier**

✓ even though the *activation function* is non-linear

It is a **discriminative model**

✓ models decision boundary directly, $P(y | \mathbf{x})$ in this case

Conditional probabilities



$$P(a|b) = \frac{P(a,b)}{P(b)}$$

$P(X, Y)$

X	Y	P
+x	+y	0.2
+x	-y	0.3
-x	+y	0.4
-x	-y	0.1

$P(+x|+y) = 0.2/0.6$ $P(-x|+y) = 0.4/0.6$ $P(-y|+x) = 0.3/0.5$

<https://inst.eecs.berkeley.edu/~cs188/fa19/assets/slides/lec13.pdf>

Set up

$$\mathcal{D} = \{(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(n)}, y^{(n)})\}$$

$$\mathbf{x}^{(i)} \in \mathbb{R}^d$$

$$y^{(i)} \in \{-1, +1\}$$

Probabilistic interpretation

$$h(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}} = \frac{e^{\mathbf{w}^T \mathbf{x}}}{e^{\mathbf{w}^T \mathbf{x}} + 1}$$

(probability of class +1) $P(y = +1 | \mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}} = \sigma(\mathbf{w}^T \mathbf{x})$

$$P(y = -1 | \mathbf{x}) = 1 - P(y = +1 | \mathbf{x})$$

(probability of class -1) $P(y = -1 | \mathbf{x}) = \frac{1}{1 + e^{\mathbf{w}^T \mathbf{x}}} = \sigma(-\mathbf{w}^T \mathbf{x})$

Probabilistic interpretation

(probability of class +1) $P(y = +1 | \mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}} = \sigma(\mathbf{w}^T \mathbf{x})$

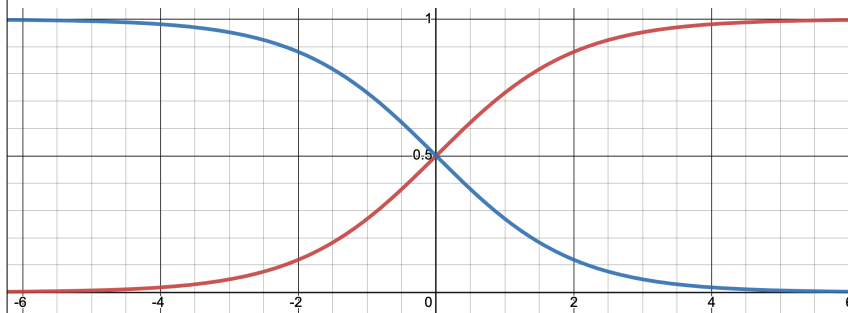
(probability of class -1) $P(y = -1 | \mathbf{x}) = \frac{1}{1 + e^{\mathbf{w}^T \mathbf{x}}} = \sigma(-\mathbf{w}^T \mathbf{x})$

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

Decision boundary

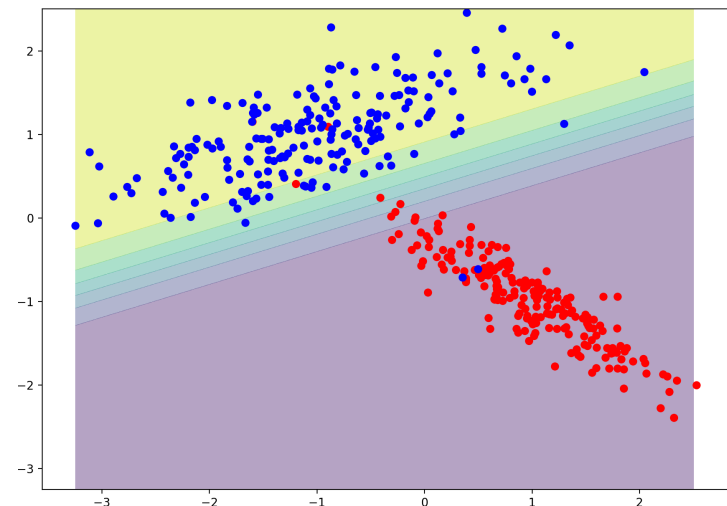
$$P(y = +1 | \mathbf{x}) = P(y = -1 | \mathbf{x}) = 0.5$$

$\sigma(\mathbf{w}^T \mathbf{x}) \qquad \sigma(-\mathbf{w}^T \mathbf{x})$

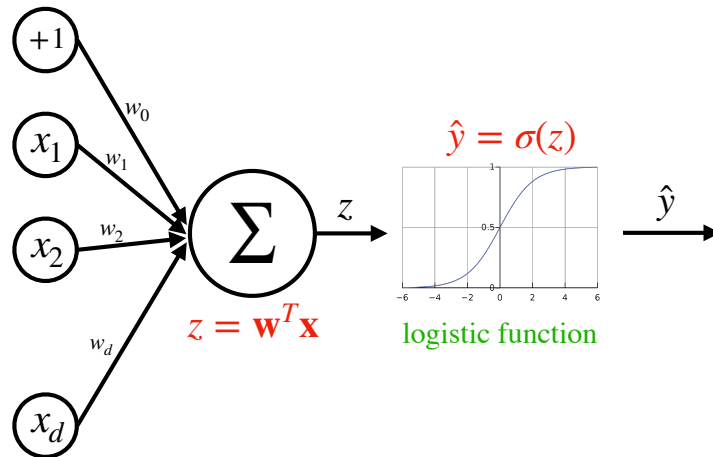


Logistic regression has a linear decision boundary $\mathbf{w}^T \mathbf{x} = 0$

Linear decision boundary



Logistic regression



Solving logistic regression

Maximum likelihood estimation (MLE)

- ▶ MLE estimates **parameters** based on the principle that if we observe \mathcal{D} , we should choose the parameters that make \mathcal{D} most probable
- ▶ We can derive formulas for W that maximize $p(\mathcal{D}; W)$
 - ✓ many machine learning algorithms follow this **maximum likelihood** principle
 - ✓ want $P(y \mid \mathbf{x}; W)$
 - ✓ learn $W^* = \arg \max_{\mathbf{w}} P(y \mid \mathbf{x}; W)$

Maximum likelihood estimation (MLE)

- ▶ The **conditional data likelihood** $\mathcal{L}(\mathbf{w})$ is the probability of the observed labels y conditioned on the feature values \mathbf{x}

✓ weights can be learned by maximizing this likelihood

$$\mathcal{L}(\mathbf{w}) = P(y^{(1)}, \dots, y^{(n)} \mid \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}; \mathbf{w}) = \prod_{i=1}^n P(y^{(i)} \mid \mathbf{x}^{(i)}; \mathbf{w})$$

$$\mathbf{w}^* = \arg \max_{\mathbf{w}} \prod_{i=1}^n P(y^{(i)} \mid \mathbf{x}^{(i)}; \mathbf{w})$$

Maximum likelihood estimation

$$\begin{aligned}
 \mathbf{w}^* &= \arg \max_{\mathbf{w}} \prod_{i=1}^n P(y^{(i)} | \mathbf{x}^{(i)}; \mathbf{w}) \\
 &= \arg \max_{\mathbf{w}} \log \left(\prod_{i=1}^n P(y^{(i)} | \mathbf{x}^{(i)}; \mathbf{w}) \right) \\
 &= \arg \max_{\mathbf{w}} \frac{1}{n} \log \left(\prod_{i=1}^n P(y^{(i)} | \mathbf{x}^{(i)}; \mathbf{w}) \right) \frac{1}{1 + e^{-y^{(i)} \mathbf{w}^T \mathbf{x}^{(i)}}} \\
 &= \arg \max_{\mathbf{w}} -\frac{1}{n} \sum_{i=1}^n \log \left(1 + e^{-y^{(i)} \mathbf{w}^T \mathbf{x}^{(i)}} \right)
 \end{aligned}$$

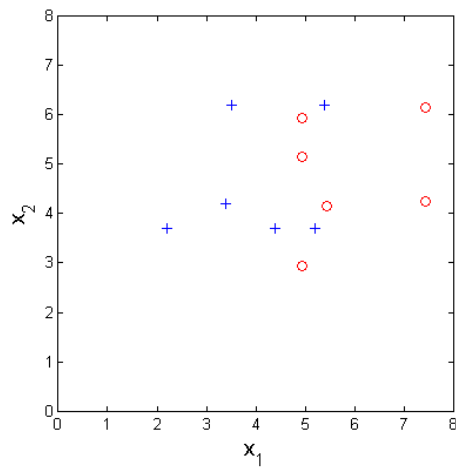
negative log
likelihood = $\arg \min_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n \log \left(1 + e^{-y^{(i)} \mathbf{w}^T \mathbf{x}^{(i)}} \right)$ error (loss)
 $e(h(\mathbf{x}), y)$

Loss function

$$L(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \log \left(1 + e^{-y^{(i)} \mathbf{w}^T \mathbf{x}^{(i)}} \right) \text{ cross-entropy loss (over a dataset)}$$

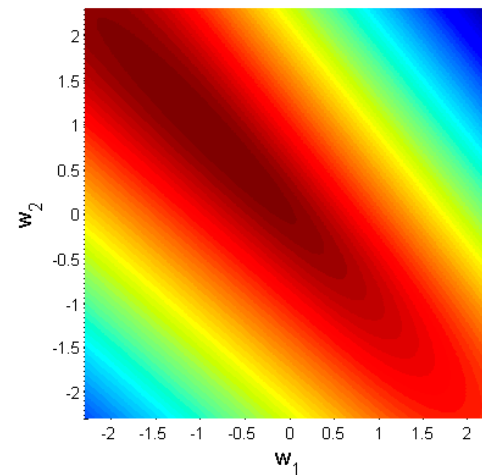
- ▶ no closed-form solution (non-linear function), but loss is convex
- ▶ can use gradient descent or second-order methods

Example: 2d dataset



<https://alliance.seas.upenn.edu/~cis520/wiki/index.php?n=Lectures.Logistic>

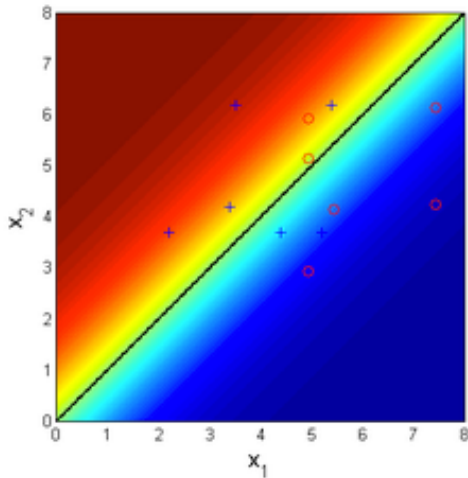
Example: loss function



plot shows contour
lines in the space
of parameters w_1
and w_2 ,
 w_0 is omitted

<https://alliance.seas.upenn.edu/~cis520/wiki/index.php?n=Lectures.Logistic>

Solution



<https://alliance.seas.upenn.edu/~cis520/wiki/index.php?n=Lectures.Logistic>

Logistic function (derivative)

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

$$\begin{aligned}\sigma'(x) &= \frac{(1 + e^{-x})(0) - (1)(-e^{-x})}{(1 + e^{-x})^2} \\ &= \frac{e^{-x}}{(1 + e^{-x})^2} \\ &= \sigma(x)(1 - \sigma(x))\end{aligned}$$

Gradient

$$\nabla_{\mathbf{w}} L(\mathbf{w}) = \left[\frac{\partial L(\mathbf{w})}{\partial w_0}, \dots, \frac{\partial L(\mathbf{w})}{\partial w_d} \right]$$

$$\begin{aligned}\frac{\partial L(\mathbf{w})}{\partial w_j} &= \frac{\partial}{\partial w_j} \frac{1}{n} \sum_{i=1}^n \log \left(1 + e^{-y^{(i)} \mathbf{w}^T \mathbf{x}^{(i)}} \right) \\ &= -\frac{1}{n} \sum_{i=1}^n \sigma \left(-y^{(i)} \mathbf{w}^T \mathbf{x}^{(i)} \right) y^{(i)} x_j^{(i)}\end{aligned}$$

How to classify new data?

► Once the final hypothesis $h(\mathbf{x})$ is known ...

✓ $h(\mathbf{x}) = p(+1 | \mathbf{x})$

✓ predict label **+1** to input instance \mathbf{x}

✓ if $p(+1 | \mathbf{x}) \geq 0.5$

✓ predict label **-1** to input instance \mathbf{x}

✓ if $p(+1 | \mathbf{x}) < 0.5$

Final remarks

- Simple classifier with **probabilistic outputs**
- Loss function is convex and can be trained with GD methods (**no closed-form**)
- Robust to overfitting
- Offers interpretability to weights (feature importance)
- However, **decision boundary is still linear**

Colab notebook

<https://colab.research.google.com/drive/1In0fRIDMQGzVQoOwTHrcyTF8G1xtmeBU?usp=sharing>