

## Binary Classification

$f($    $) \rightarrow \{\text{Cat, Dog}\}$

$y = f(\mathbf{x})$ ,   Input:  $\mathbf{x} \in \mathbb{R}^n \rightarrow$  Output:  $y \in \{0, 1\}$

Linear Reg.

$\mathbb{R}$

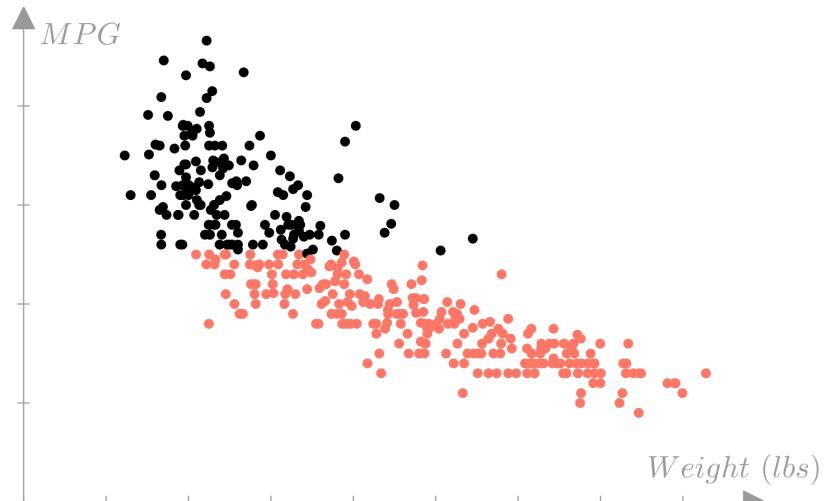
# Binary Classification

$$y = f(\mathbf{x}), \quad \text{Input: } \mathbf{x} \in \mathbb{R}^n \longrightarrow \text{Output: } y \in \{0, 1\}$$

Input:  $\mathbf{x}_i = \begin{bmatrix} \text{Weight} \\ \text{Horsepower} \\ \text{Displacement} \\ \text{0-60mph} \end{bmatrix}$ , Output:  $y_i = \begin{cases} 1 : \text{Meets target } (MPG \geq 30) \\ 0 : \text{Fails to meet target } (MPG < 30) \end{cases}$

Honda Accord:  $\begin{bmatrix} \text{Weight:} & 2500 \text{ lbs} \\ \text{Horsepower:} & 123 \text{ HP} \\ \text{Displacement:} & 2.4 \text{ L} \\ \text{0-60mph:} & 7.8 \text{ Sec} \end{bmatrix} \rightarrow 1 \text{ (Meets target)}$

Dodge Aspen:  $\begin{bmatrix} \text{Weight:} & 3800 \text{ lbs} \\ \text{Horsepower:} & 155 \text{ HP} \\ \text{Displacement:} & 3.2 \text{ L} \\ \text{0-60mph:} & 6.8 \text{ Sec} \end{bmatrix} \rightarrow 0 \text{ (Does not meet target)}$



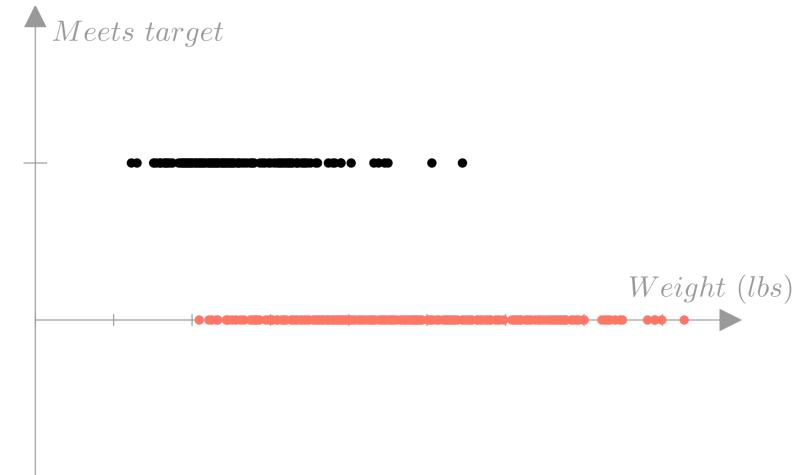
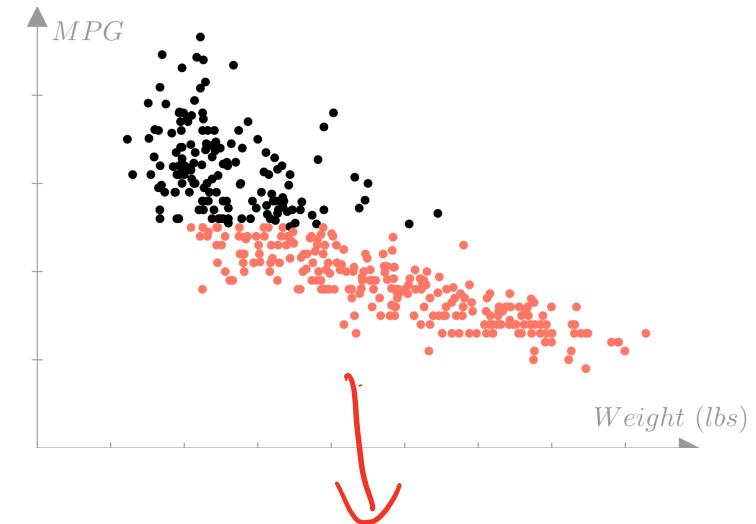
# Binary Classification

$$y = f(\mathbf{x}), \quad \text{Input: } \mathbf{x} \in \mathbb{R}^n \longrightarrow \text{Output: } y \in \{0, 1\}$$

$$\text{Input: } \mathbf{x}_i = \begin{bmatrix} \text{Weight} \\ \text{Horsepower} \\ \text{Displacement} \\ \text{0-60mph} \end{bmatrix}, \quad \text{Output: } y_i = \begin{cases} 1 : \text{Meets target } (MPG \geq 30) \\ 0 : \text{Fails to meet target } (MPG < 30) \end{cases}$$

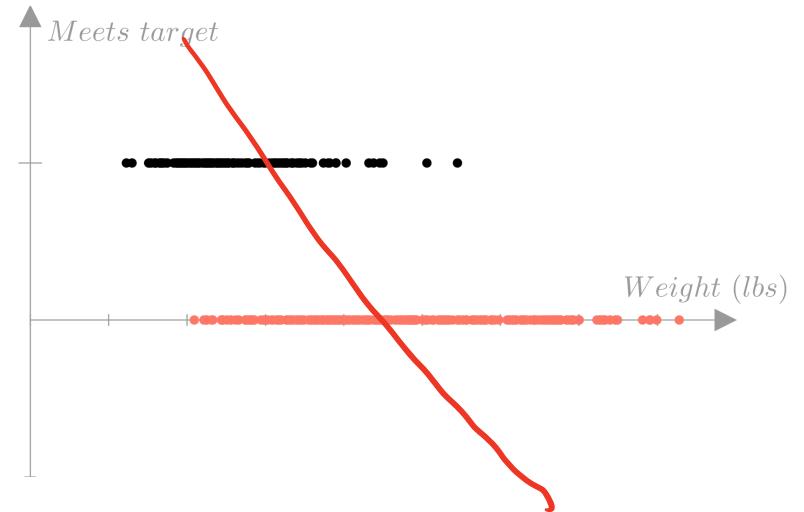
Honda Accord:  $\begin{bmatrix} \text{Weight:} & 2500 \text{ lbs} \\ \text{Horsepower:} & 123 \text{ HP} \\ \text{Displacement:} & 2.4 \text{ L} \\ \text{0-60mph:} & 7.8 \text{ Sec} \end{bmatrix} \rightarrow 1 \text{ (Meets target)}$

Dodge Aspen:  $\begin{bmatrix} \text{Weight:} & 3800 \text{ lbs} \\ \text{Horsepower:} & 155 \text{ HP} \\ \text{Displacement:} & 3.2 \text{ L} \\ \text{0-60mph:} & 6.8 \text{ Sec} \end{bmatrix} \rightarrow 0 \text{ (Does not meet target)}$



# Can we do binary classification with linear regression?

$$f(\mathbf{x}) = \mathbf{x}^T \mathbf{w} + b = \sum_{i=1}^d x_i w_i + b$$



$$\mathbf{w} = \begin{bmatrix} w_1 \\ w_2 \\ w_3 \\ b \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ 1 \end{bmatrix}$$

$$f(\mathbf{x}) = \mathbf{x}^T \mathbf{w} \rightarrow f(\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbf{x}^T \mathbf{w} \geq 0 \\ 0 & \text{if } \mathbf{x}^T \mathbf{w} < 0 \end{cases}$$

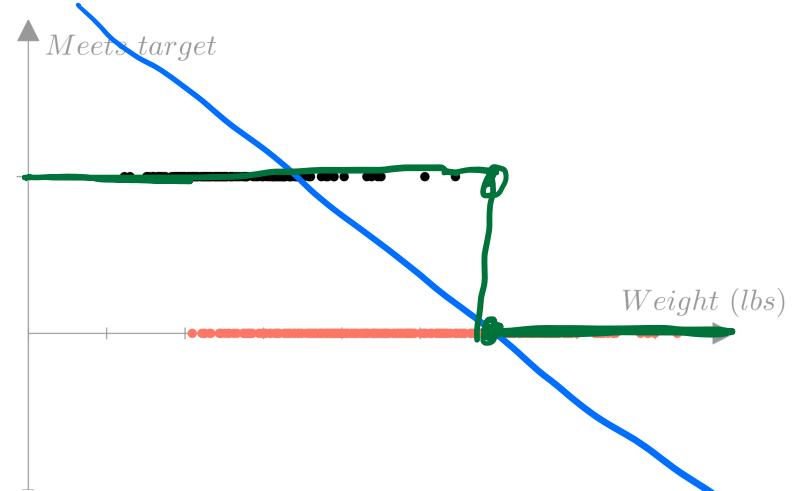
Linear fcs.

$$f(\mathbf{x}) = \mathbb{I}(\mathbf{x}^T \mathbf{w} \geq 0)$$

Indicator func.

Binary predictions with thresholding

$$\mathbf{x}^T \mathbf{w} \equiv \sum_{i=1}^d x_i w_i + b$$



## Binary predictions with thresholding

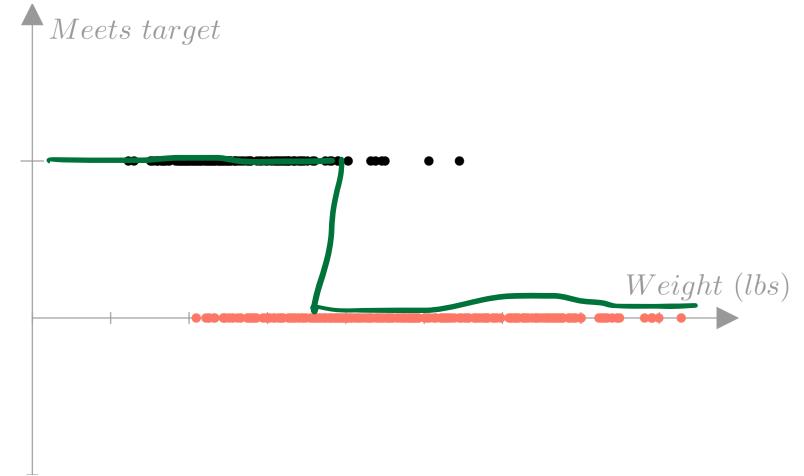
$$f(\mathbf{x}) = \mathbf{x}^T \mathbf{w} \quad \longrightarrow \quad f(\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbf{x}^T \mathbf{w} \geq 0 \\ 0 & \text{if } \mathbf{x}^T \mathbf{w} < 0 \end{cases}$$

$$f(\mathbf{x}) = \mathbb{I}(\mathbf{x}^T \mathbf{w} \geq 0)$$

Meets target =  $f(\mathbf{x}) =$

$$((\text{weight})w_1 + (\text{horsepower})w_2 + (\text{displacement})w_3 + (0\text{-}60\text{mph})w_4 + b) \geq 0$$

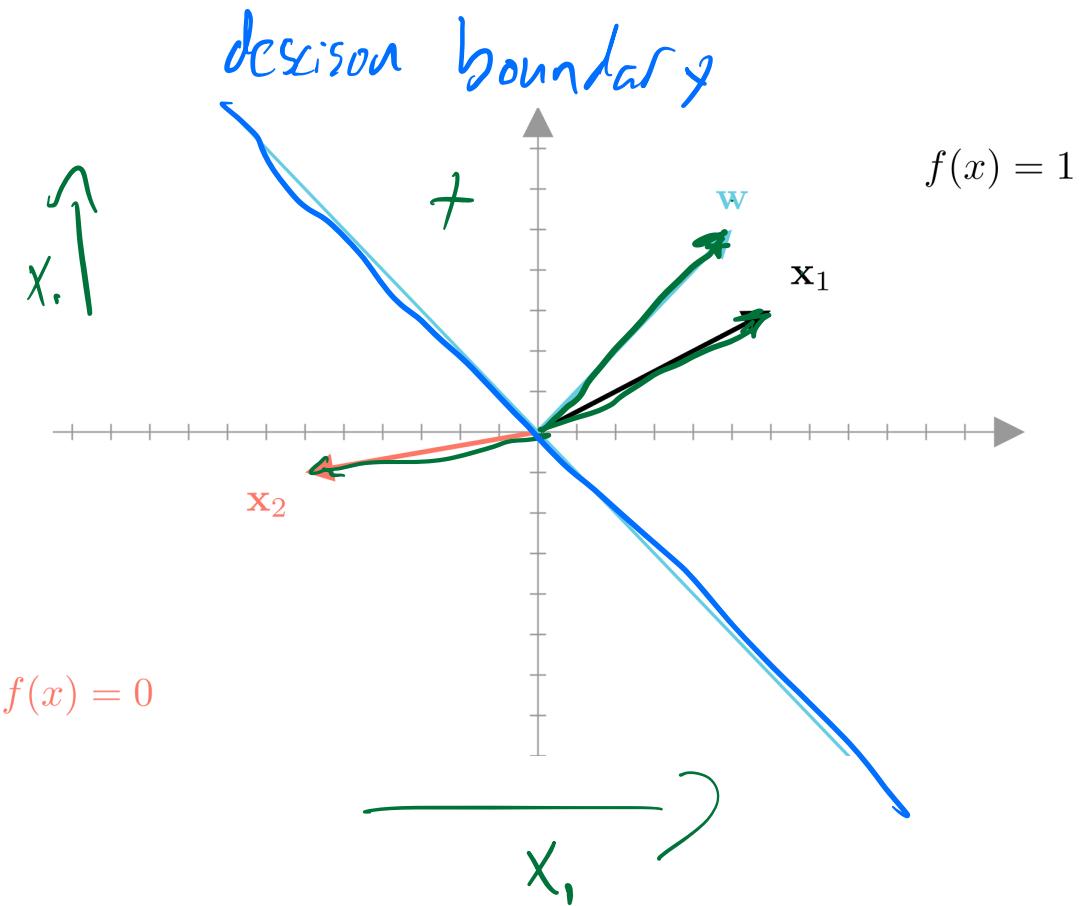
$$f(\mathbf{x}) = \left( \begin{bmatrix} \text{Weight} \\ \text{Horsepower} \\ \text{Displacement} \\ 0\text{-}60\text{mph} \\ 1 \end{bmatrix} \cdot \begin{bmatrix} w_1 \\ w_2 \\ w_3 \\ w_4 \\ b \end{bmatrix} \geq 0 \right)$$



# Thresholds in higher dimensions

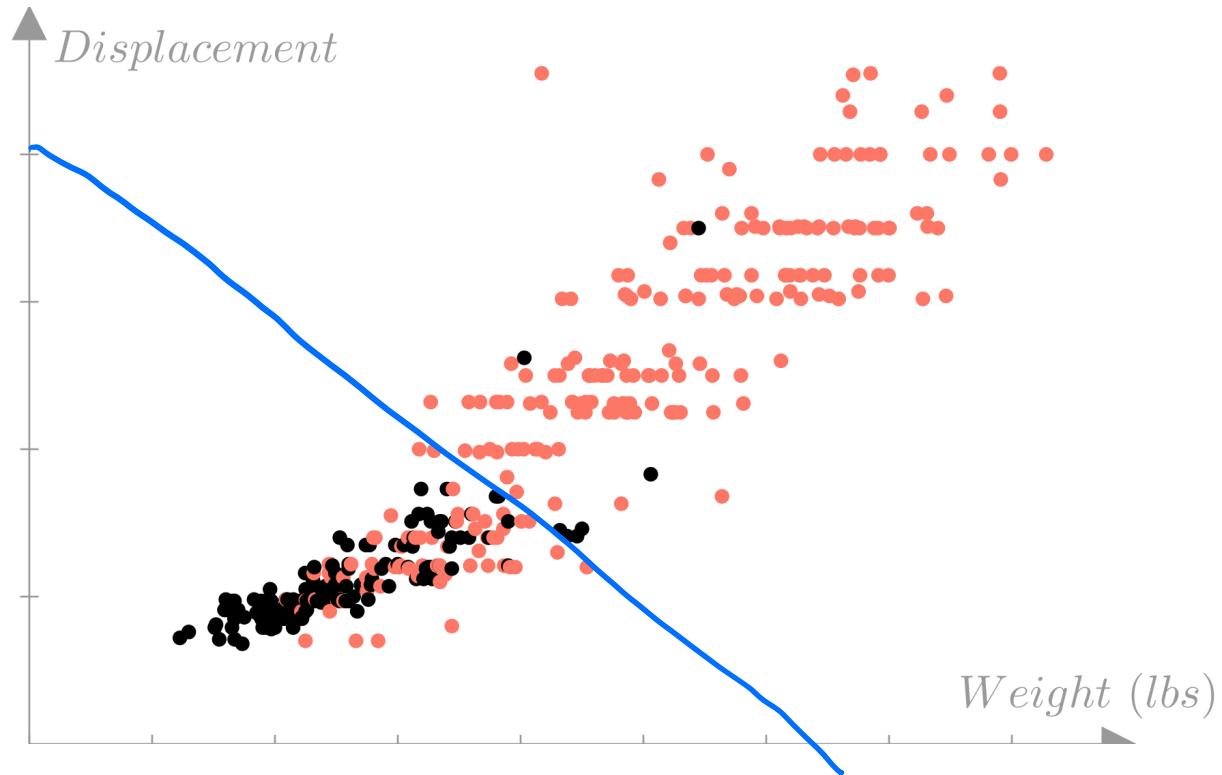
$$f(\mathbf{x}) = \mathbf{x}^T \mathbf{w} \rightarrow f(\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbf{x}^T \mathbf{w} \geq 0 \\ 0 & \text{if } \mathbf{x}^T \mathbf{w} < 0 \end{cases}$$

$$\mathbf{x}^T \mathbf{w} = \|\mathbf{x}\|_2 \|\mathbf{w}\|_2 \cos \theta$$



## Decision boundaries

$$f(\mathbf{x}) = \mathbf{x}^T \mathbf{w} \quad \rightarrow \quad f(\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbf{x}^T \mathbf{w} \geq 0 \\ 0 & \text{if } \mathbf{x}^T \mathbf{w} < 0 \end{cases}$$



# Accuracy as a loss

Accuracy:

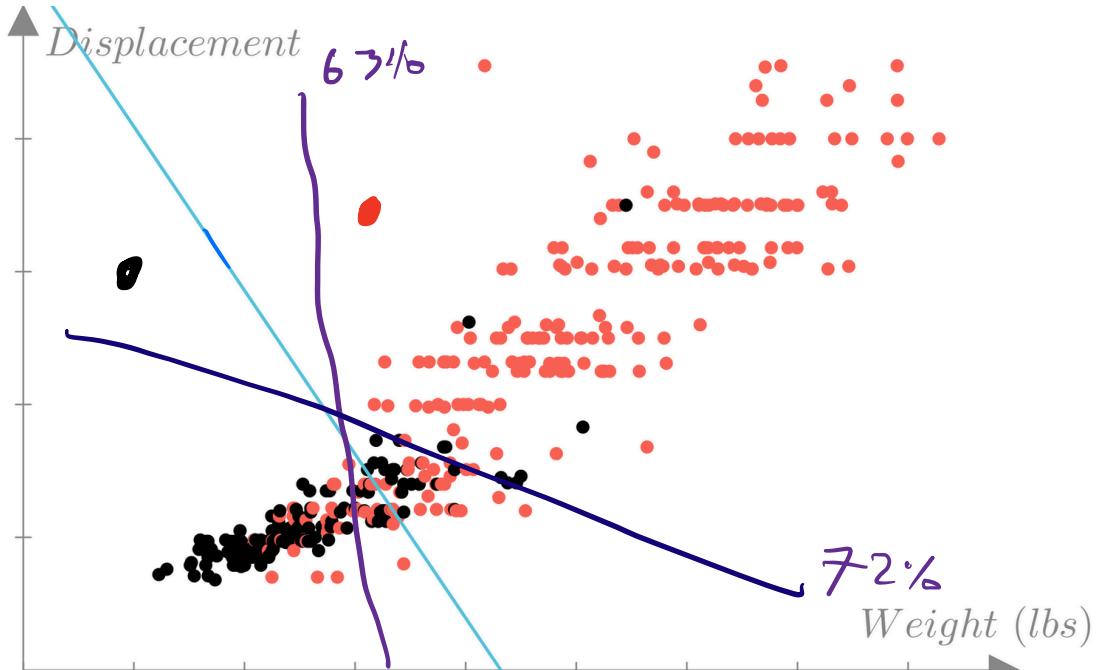
$$\frac{\text{\# of correct predictions}}{\text{Total predictions}}$$

Prediction function  
 $\mathbb{I}(x^T w \geq 0)$

$$\text{Accuracy} = \frac{1}{N} \sum_{i=1}^N \mathbb{I}(f(\mathbf{x}_i) = y_i)$$

Dataset      |      label  
1            0

Accuracy: 0.8291



Can we use accuracy to find optimal parameters?

Gradient descent

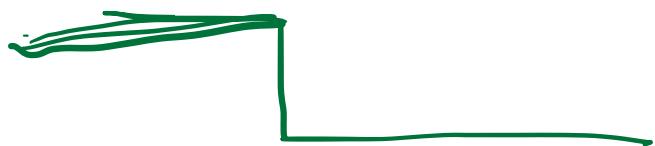
While  $\nabla f(\mathbf{w}^{(i)}) \neq \mathbf{0}$ :  $\mathbf{w}^{(i+1)} \leftarrow \mathbf{w}^{(i)} - \nabla f(\mathbf{w}^{(i)})$

Loss

$$\text{Accuracy} = \frac{1}{N} \sum_{i=1}^N \mathbb{I}(f(\mathbf{x}_i) = y_i)$$

$$\frac{\partial}{\partial w} \frac{1}{N} \sum_{i=1}^N \mathbb{I}(f(x_i) = y_i)$$

$$\frac{\partial}{\partial w} \mathbb{I}(f_w(x_i) = y_i)$$



$P(y|x, \omega)$

## The Bernoulli distribution

Probability of heads:  $q$ , Probability of tails:  $1 - q$

coin flip

PMF

$$p(y) = \begin{cases} q & \text{if } y = 1 \\ 1 - q & \text{if } y = 0 \end{cases} \quad \begin{matrix} \underline{q \in [0, 1]}, \\ \underline{y \in \{0, 1\}} \end{matrix}$$

# The Bernoulli distribution

Probability of **heads**:  $q$ , Probability of **tails**:  $1 - q$

$$p(y) = \begin{cases} q & \text{if } y = 1 \\ 1 - q & \text{if } y = 0 \end{cases} \quad q \in [0, 1], y \in \{0, 1\}$$

$$p(y) = q^y (1 - q)^{1-y}$$

Annotations:

- A red bracket above the term  $(1 - q)^{1-y}$  is labeled  $y=0$ .
- A red bracket above the term  $q^y$  is labeled  $y=1$ .
- A purple bracket under the entire expression  $q^y (1 - q)^{1-y}$  is labeled  $y$ .

$$\log p(y) = y \log q + (1 - y) \log(1 - q)$$

Annotation:

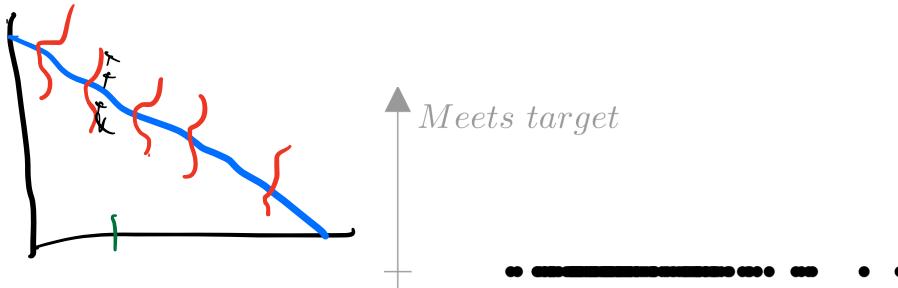
- A purple bracket under the terms  $y \log q$  and  $(1 - y) \log(1 - q)$  is labeled  $y$ .

# A probabilistic model for binary classification

Linear regression: Normal observations

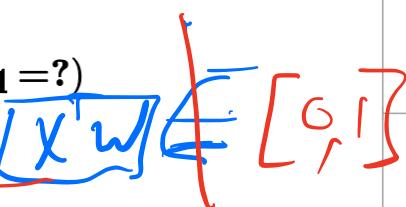
$$y_i \sim \mathcal{N}(\mathbf{x}_i^T \mathbf{w}, \sigma^2)$$

*mean*



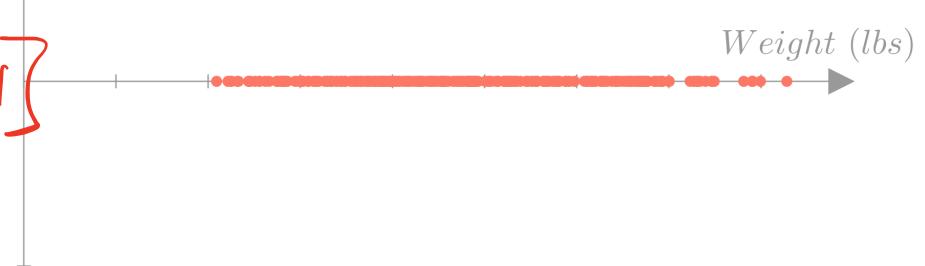
Bernoulli observations?

$$\mathbf{x}^T \mathbf{w} \notin [0, 1] \rightarrow y_i \sim \text{Bernoulli}(q = ?)$$



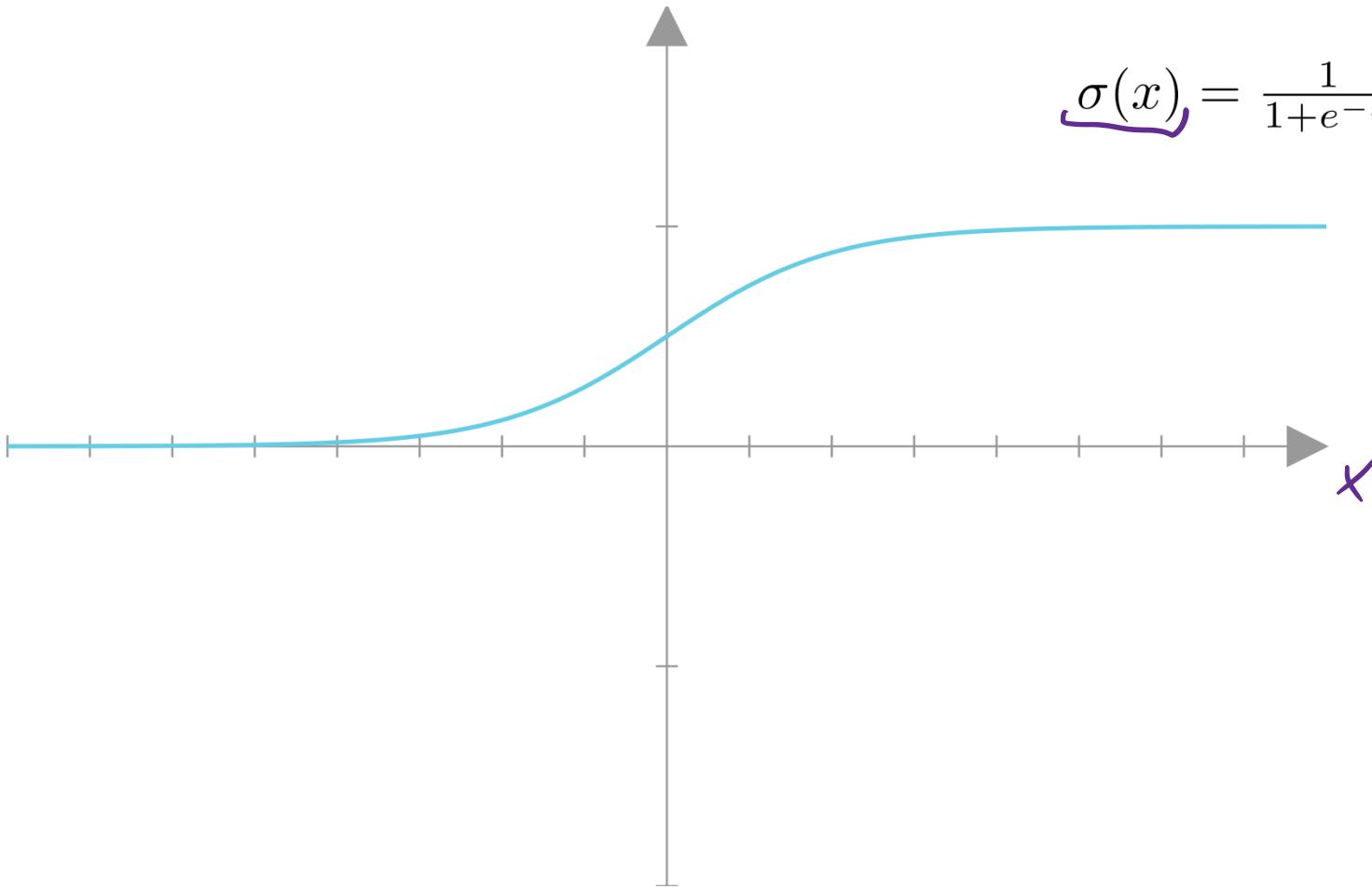
Need  $g(x)$ :  $\mathbb{R} \rightarrow [0, 1]$

**Input:**  $x \in \mathbb{R} \rightarrow$  **Output:**  $y \in [0, 1]$



## The sigmoid function

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



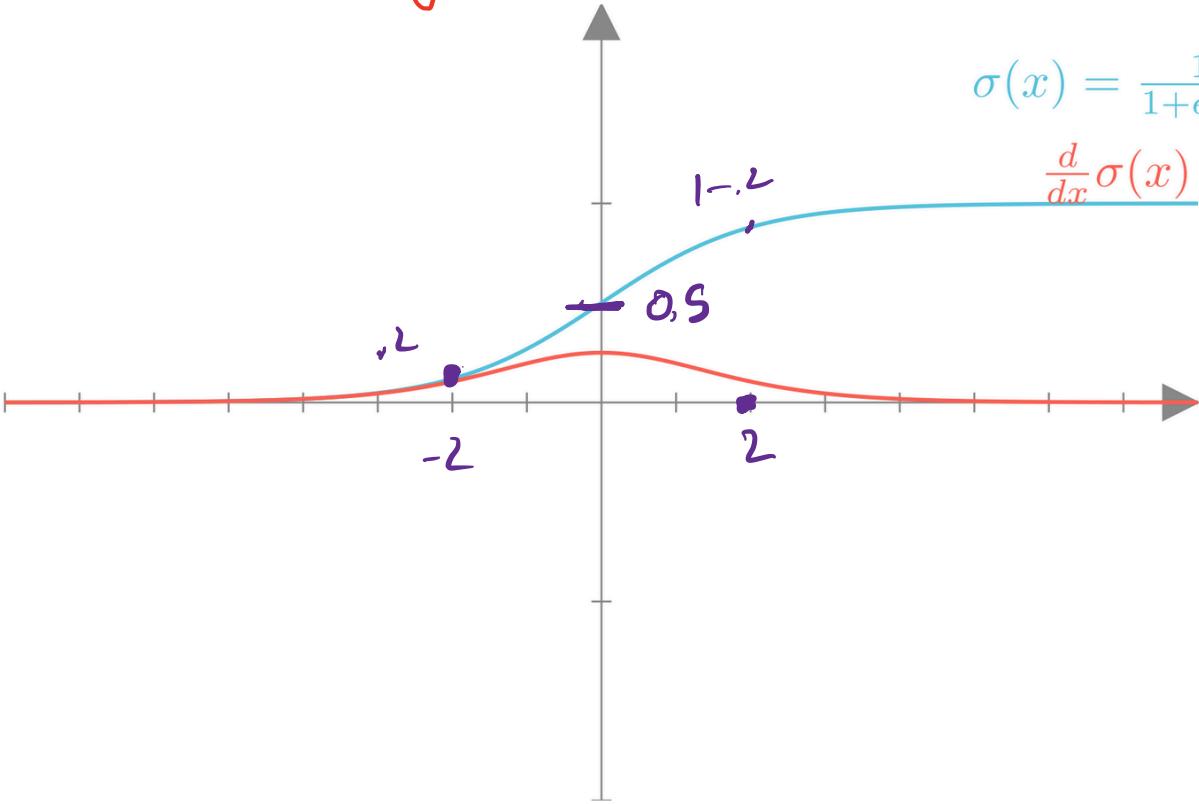
# Properties of the sigmoid function

$$\frac{d}{dx} \sigma(x) = \sigma(x)(1 - \sigma(x))$$

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

$$\sigma(0) = 0.5$$

$$1 - \sigma(x) = \sigma(-x)$$



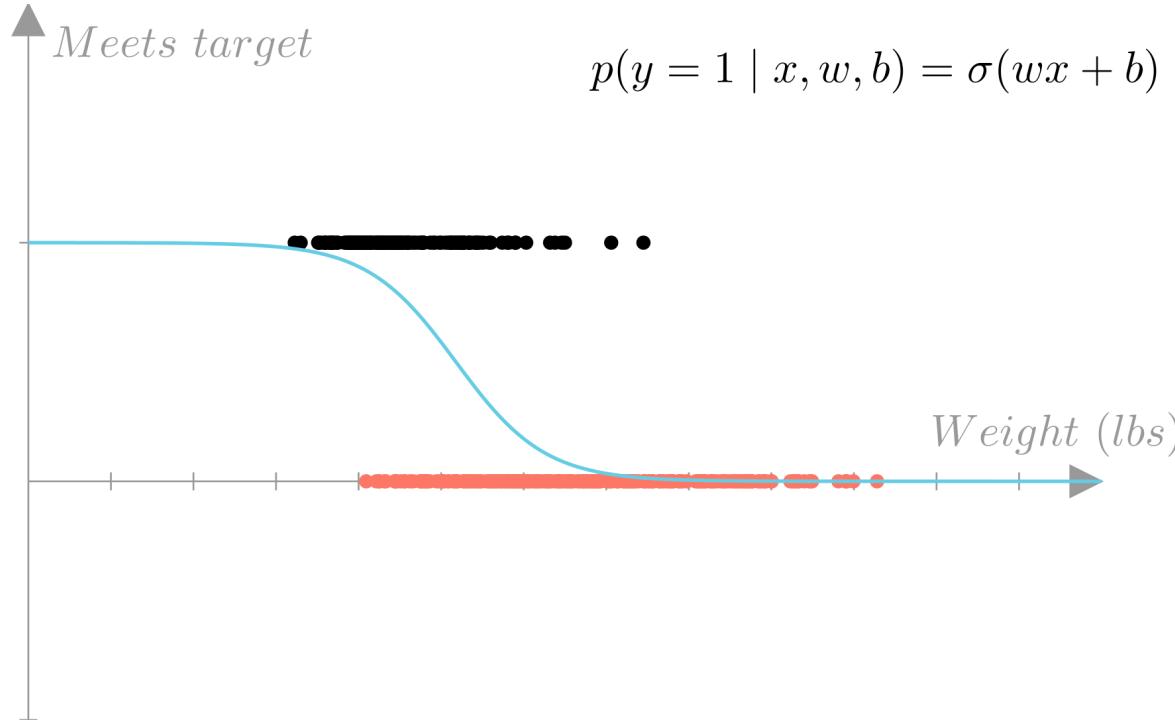
# Logistic regression

Parameters  
/

$$y_i \sim \text{Bernoulli}(\sigma(\mathbf{x}_i^T \mathbf{w}))$$

$p(y_i = 1 \mid \mathbf{x}_i, \mathbf{w}) = \sigma(\mathbf{x}_i^T \mathbf{w}),$

$$p(y_i = 0 \mid \mathbf{x}_i, \mathbf{w}) = 1 - \sigma(\mathbf{x}_i^T \mathbf{w}) = \sigma(-\mathbf{x}_i^T \mathbf{w})$$



**Prediction function:**  $f(\mathbf{x}) = \begin{cases} 1 & \text{if } p(y = 1 \mid \mathbf{x}, \mathbf{w}) \geq p(y = 0 \mid \mathbf{x}, \mathbf{w}) \\ 0 & \text{otherwise} \end{cases}$

# Logistic regression

Prediction function:  $f(\mathbf{x}) = \begin{cases} 1 & \text{if } p(y = 1 | \mathbf{x}, \mathbf{w}) \geq p(y = 0 | \mathbf{x}, \mathbf{w}) \\ 0 & \text{otherwise} \end{cases}$

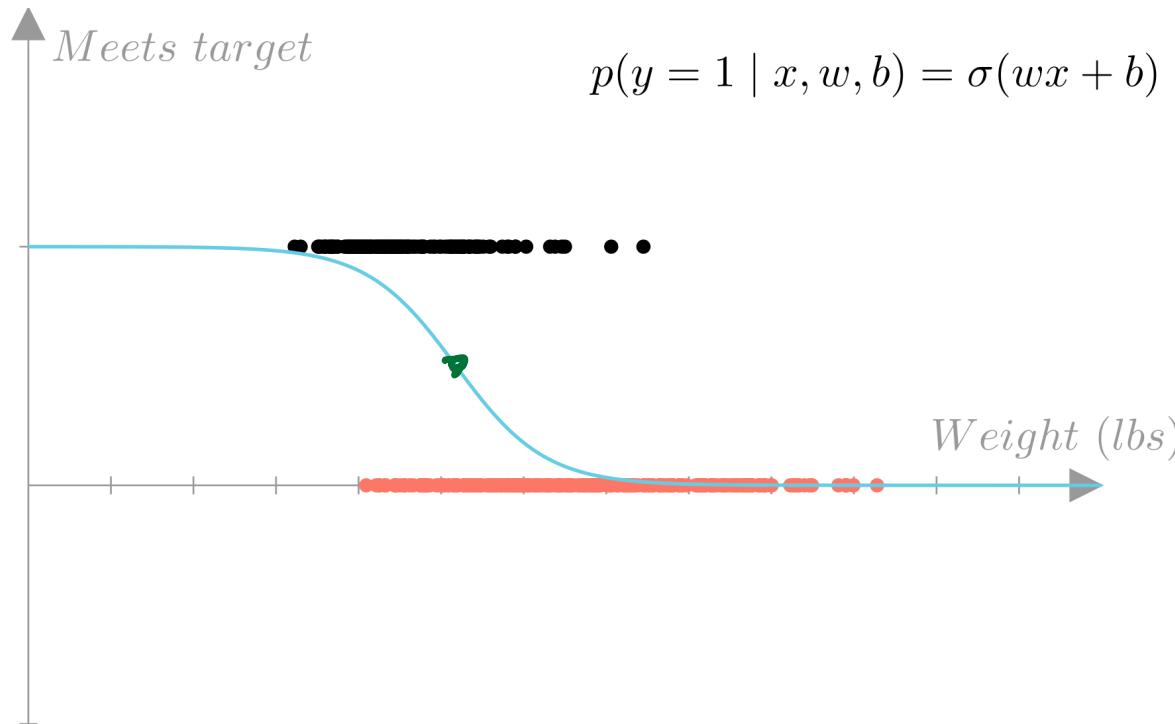
$$P(y=1|\mathbf{x}, \mathbf{w}) > P(y=0|\mathbf{x}, \mathbf{w})$$

$$\equiv P(y=1|\mathbf{x}, \mathbf{w}) \geq 0.5$$

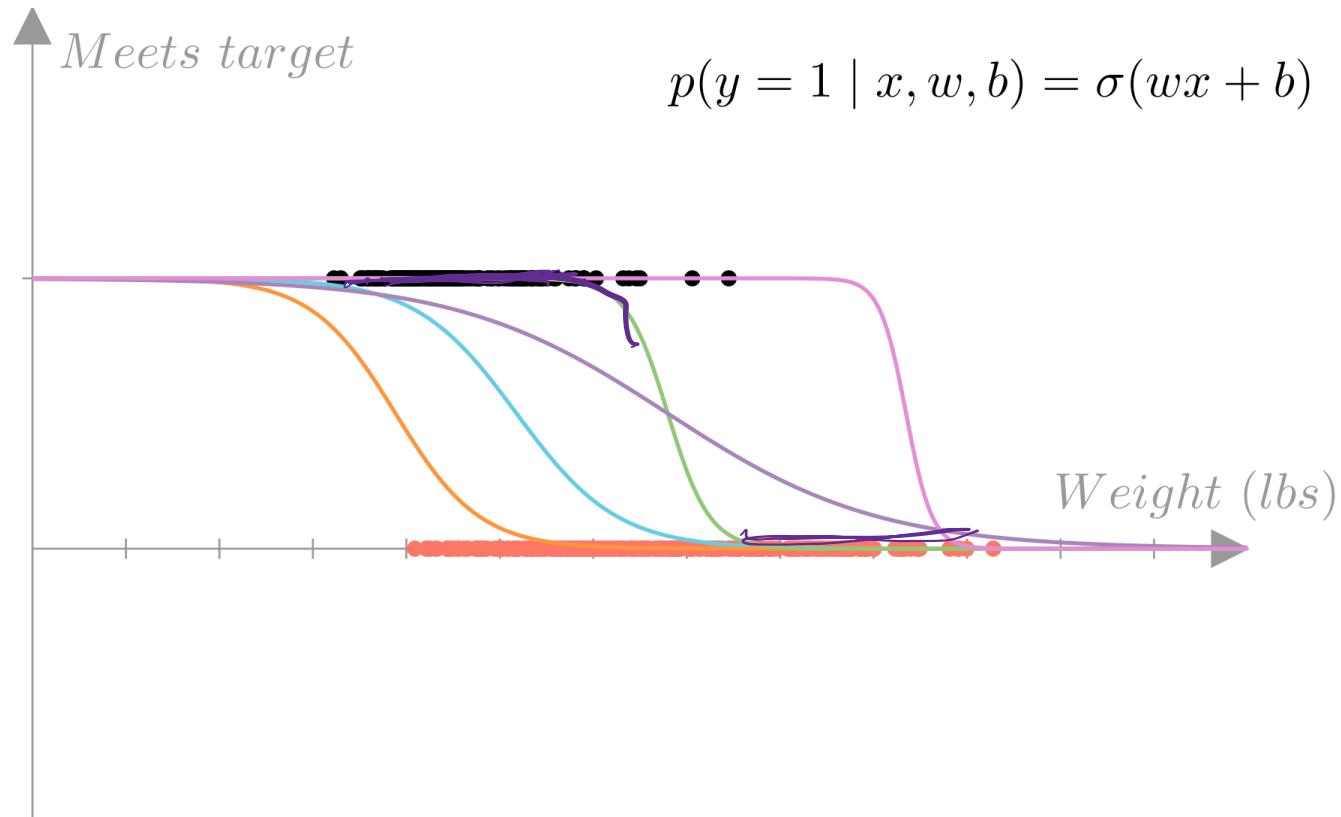
$$p(y = 1 | \mathbf{x}, \mathbf{w}) = \sigma(\mathbf{x}^T \mathbf{w}) \geq 0.5$$

$$\sigma(0) = 0.5$$

$$p(y_i = 1) \geq 0.5 \quad \rightarrow \quad \mathbf{x}^T \mathbf{w} \geq 0$$



# Maximum likelihood for logistic regression



$$\boxed{\mathbf{w}^*} = \operatorname{argmax}_{\mathbf{w}} p(\mathbf{y} \mid \mathbf{X}, \mathbf{w}) = \operatorname{argmax}_{\mathbf{w}} p(y_1, \dots, y_N \mid \mathbf{x}_1, \dots, \mathbf{x}_N, \mathbf{w})$$

# Maximum likelihood for logistic regression

$$\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmax}} p(\mathbf{y} \mid \mathbf{X}, \mathbf{w}) = \underset{\mathbf{w}}{\operatorname{argmax}} p(y_1, \dots, y_N \mid \mathbf{x}_1, \dots, \mathbf{x}_N, \mathbf{w})$$

$$\equiv \prod_{i=1}^N p(y_i \mid \mathbf{x}_i, \mathbf{w})$$

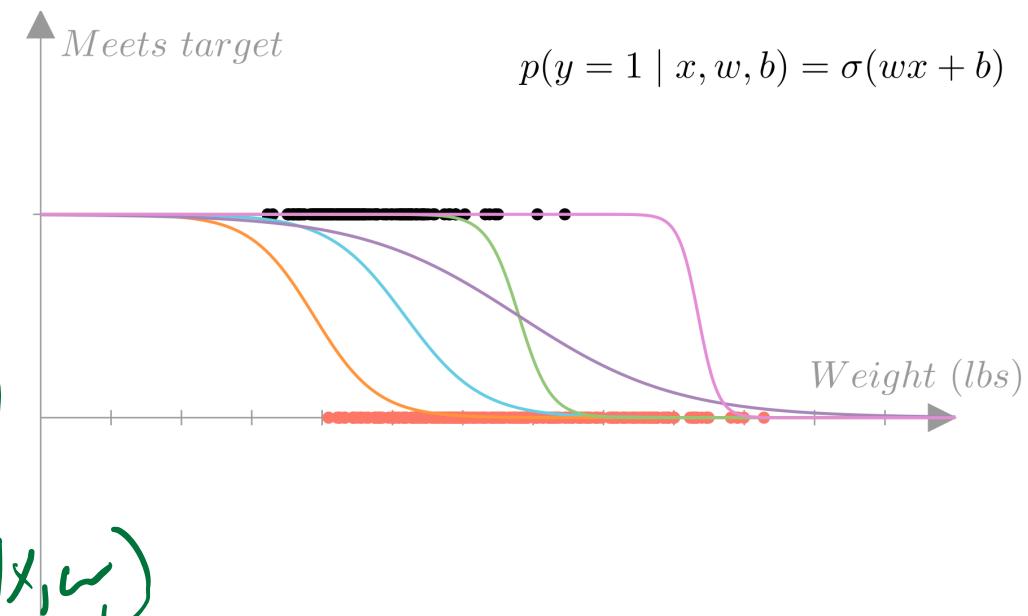
$$\text{Loss}(\mathbf{w}) = \text{NLL}(\mathbf{w}, \mathbf{X}, \mathbf{y}) = -\sum_{i=1}^N \log p(y_i \mid \mathbf{x}_i, \mathbf{w})$$

$$p(y_i = 1 \mid \mathbf{x}_i, \mathbf{w}) = \sigma(\mathbf{x}_i^T \mathbf{w}), \quad p(y_i = 0 \mid \mathbf{x}_i, \mathbf{w}) = 1 - \sigma(\mathbf{x}_i^T \mathbf{w}) = \sigma(-\mathbf{x}_i^T \mathbf{w})$$

Sum over our dataset

$$\text{NLL}(\mathbf{w}, \mathbf{X}, \mathbf{y}) = -\sum_{i=1}^N \left[ y_i \log \sigma(\mathbf{x}_i^T \mathbf{w}) + (1 - y_i) \log(1 - \sigma(\mathbf{x}_i^T \mathbf{w})) \right]$$

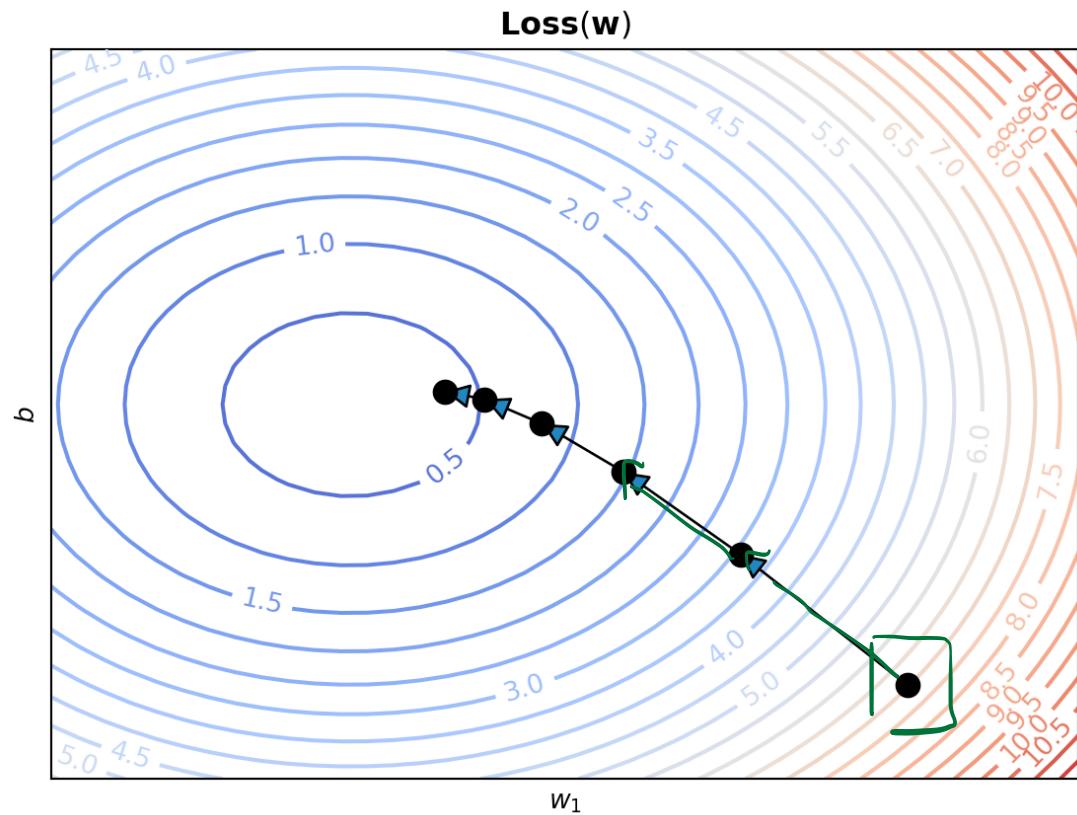
%  $\rho(y=1 \mid \mathbf{x}_i, \mathbf{w})$        $\rho(y=0 \mid \mathbf{x}_i, \mathbf{w})$



$$p(y = 1 \mid x, w, b) = \sigma(wx + b)$$

# Gradient descent

While  $\nabla f(\mathbf{w}^{(i)}) \neq \mathbf{0}$  :  $\mathbf{w}^{(i+1)} \leftarrow \mathbf{w}^{(i)} - \nabla f(\mathbf{w}^{(i)})$



## Deriving the gradient

$$\text{NLL}(\mathbf{w}, \mathbf{X}, \mathbf{y}) = - \sum_{i=1}^N \left[ y_i \log \sigma(\mathbf{x}_i^T \mathbf{w}) + (1 - y_i) \log(1 - \sigma(\mathbf{x}_i^T \mathbf{w})) \right]$$

$$\frac{\partial}{\partial \mathbf{w}} - \sum_{i=1}^N \left[ y_i \log \sigma(\mathbf{x}_i^T \mathbf{w}) + (1 - y_i) \log(1 - \sigma(\mathbf{x}_i^T \mathbf{w})) \right]$$

$$= - \sum_{i=1}^N \frac{d}{d \mathbf{w}} \left[ y_i \log \sigma(\mathbf{x}_i^T \mathbf{w}) + (1 - y_i) \log(1 - \sigma(\mathbf{x}_i^T \mathbf{w})) \right]$$

$$= - \sum_{i=1}^N \left[ y_i \frac{d}{d \mathbf{w}} \log \sigma(\mathbf{x}_i^T \mathbf{w}) + (1 - y_i) \frac{d}{d \mathbf{w}} \log(1 - \sigma(\mathbf{x}_i^T \mathbf{w})) \right]$$

$$\hookrightarrow = \frac{1}{\sigma(\mathbf{x}_i^T \mathbf{w})} \frac{d}{d \mathbf{w}} \sigma(\mathbf{x}_i^T \mathbf{w})$$

$$\frac{d}{d \mathbf{x}} \sigma(x) = \sigma(x)(1 - \sigma(x))$$

$$= \frac{1}{\sigma(\mathbf{x}_i^T \mathbf{w})} (\sigma(\mathbf{x}_i^T \mathbf{w})(1 - \sigma(\mathbf{x}_i^T \mathbf{w}))) \frac{d}{d \mathbf{w}} \mathbf{x}_i^T \mathbf{w}$$

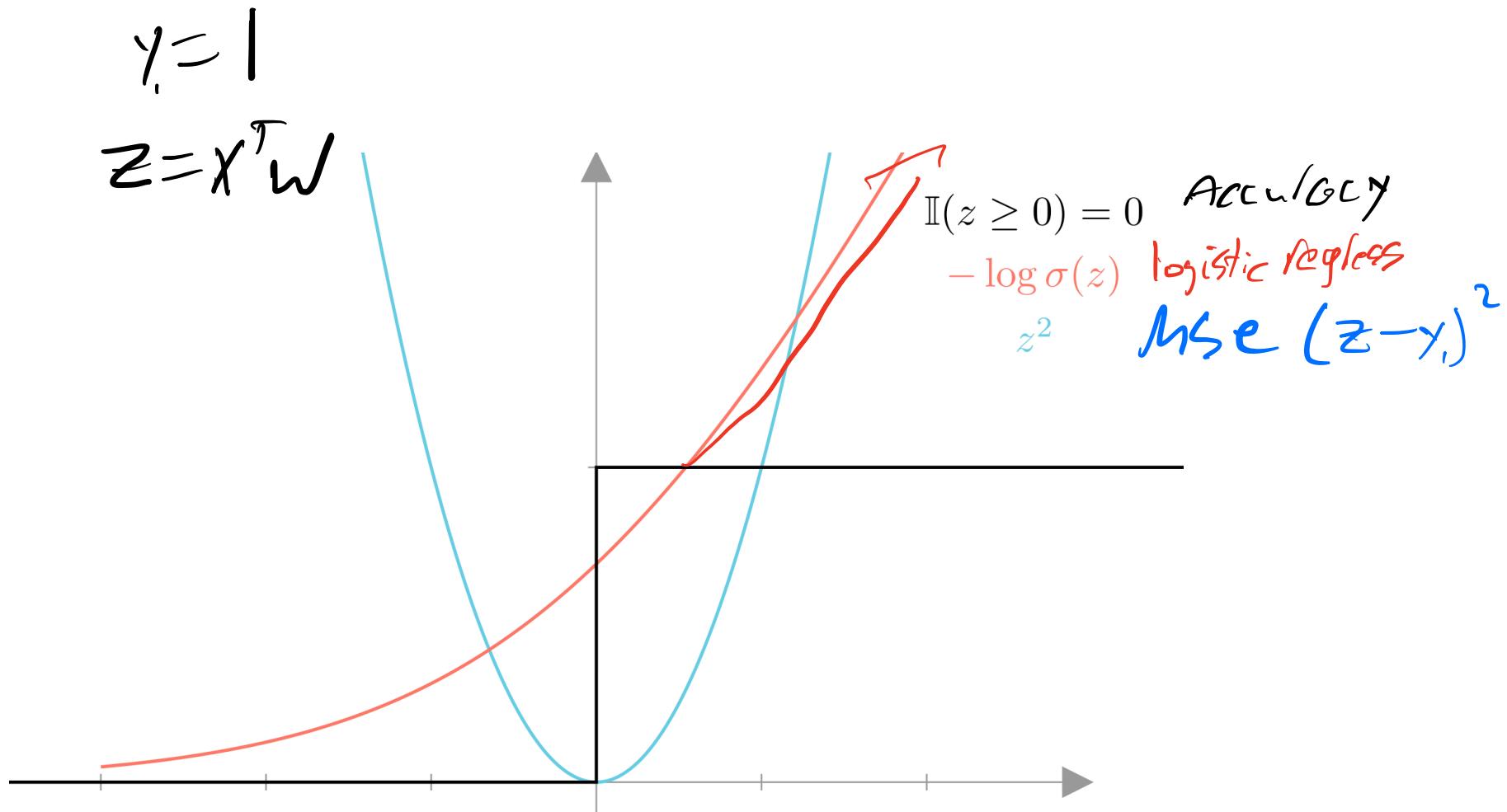
$$= \left( I - \sigma(X^T \omega) \right) X^T$$

## Deriving the gradient

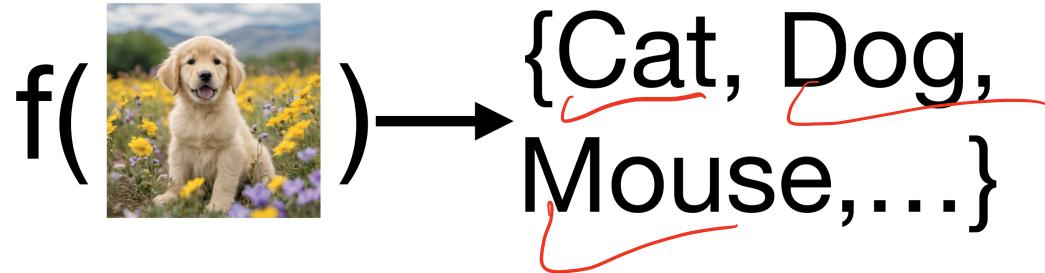
$$\text{NLL}(\mathbf{w}, \mathbf{X}, \mathbf{y}) = - \sum_{i=1}^N \left[ y_i \log \sigma(\mathbf{x}_i^T \mathbf{w}) + (1 - y_i) \log(1 - \sigma(\mathbf{x}_i^T \mathbf{w})) \right]$$

$$\nabla_{\mathbf{w}} \text{NLL}(\mathbf{w}, \mathbf{X}, \mathbf{y}) = - \sum_{i=1}^N \left( 1 - \sigma((2y_i - 1)\mathbf{x}_i^T \mathbf{w}) \right) \left( (2y_i - 1)\mathbf{x}_i \right)$$


## Comparing loss functions



## Multi-class classification



$$y = f(\mathbf{x}), \quad \text{Input: } \mathbf{x} \in \mathbb{R}^n \longrightarrow \text{Output: } y \in \{1, 2, \dots, C\}$$

Ordering irrelevant!

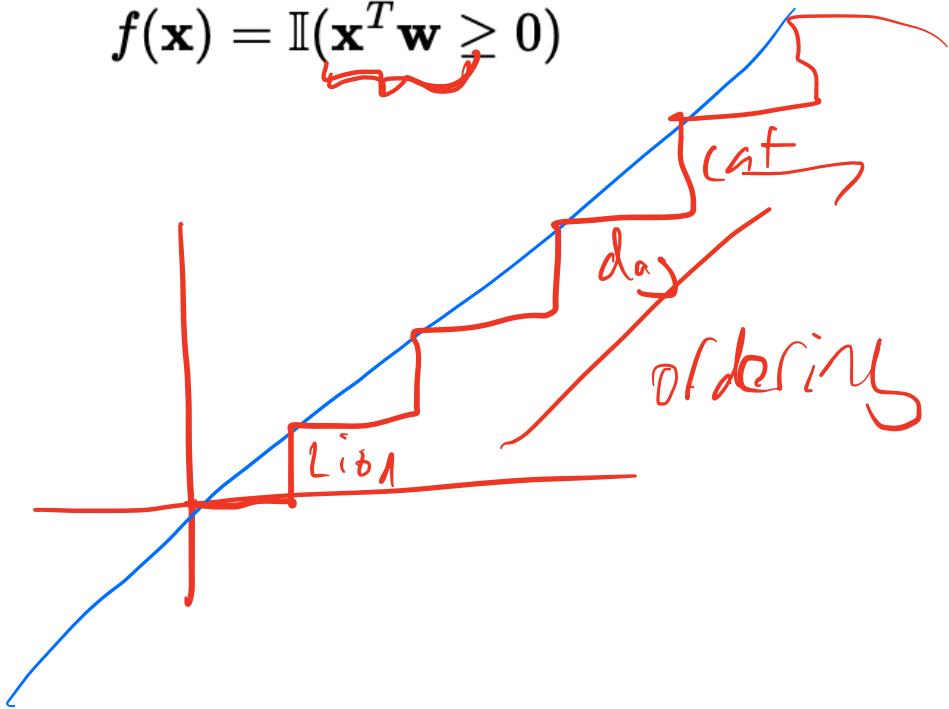
1: Cat, 2: Dog, 3: Mouse

1: Dog, 2: Mouse, 3: Cat

## Multi-class prediction functions

### Binary thresholding

$$f(\mathbf{x}) = \mathbb{I}(\mathbf{x}^T \mathbf{w} \geq 0)$$



### Multi-class thresholding

$$f(\mathbf{x}) = \operatorname{argmax}_{c \in \{1 \dots C\}} \mathbf{x}^T \mathbf{w}_c$$

↳ different function

## Multi-class prediction functions

### Multi-class thresholding

$$f(\mathbf{x}) = \operatorname{argmax}_{c \in \{1 \dots C\}} \mathbf{x}^T \mathbf{w}_c$$

For convenience, we can also define a matrix that contains all  $C$  parameter vectors:

*Capital W*

$$\mathbf{W} = \begin{bmatrix} \mathbf{w}_1^T \\ \mathbf{w}_2^T \\ \vdots \\ \mathbf{w}_C^T \end{bmatrix} = \begin{bmatrix} W_{11} & W_{12} & \dots & W_{1d} \\ W_{21} & W_{22} & \dots & W_{2d} \\ \vdots & \vdots & \ddots & \vdots \\ W_{C1} & W_{C2} & \dots & W_{Cd} \end{bmatrix}$$

*( sets of parameters  $\mathbf{w}_1, \dots, \mathbf{w}_C$ )*

With this notation, our prediction function becomes:

$$f(\mathbf{x}) = \operatorname{argmax}_{c \in \{1 \dots C\}} (\mathbf{x}^T \mathbf{W}^T)_c, \quad \mathbf{W} \in \mathbb{R}^{C \times d}$$

$$\mathbf{x}^T \mathbf{w}_c = (\mathbf{x}^T \mathbf{W}^T)_c$$

## Multi-class decision boundaries

$$f(\mathbf{x}) = \operatorname{argmax}_{c \in \{0,1\}} (\mathbf{x}^T \mathbf{W}^T)_c = \mathbb{I}(\mathbf{x}^T \mathbf{w}_1 - \mathbf{x}^T \mathbf{w}_0 \geq 0)$$

## Categorical distribution

$$p(y = c) = q_c, \quad y \in \{1 \dots C\}$$

$$\mathbf{q} \in \mathbb{R}^C \quad q_c \geq 0 \quad \forall c \in \{1 \dots C\} \quad \sum_{c=1}^C q_c = 1$$

## A probabilistic model for multi-class classification

$$y_i \sim \text{Categorical}(\mathbf{q} = ?), \quad \mathbf{q} = \mathbf{x}_i^T \mathbf{W}^T ?$$

$$\mathbf{x}^T \mathbf{W}^T \in \mathbb{R}^C, \quad q_c \geq 0 \quad \forall c \in \{1 \dots C\}, \quad \sum_{c=1}^C q_c = 1$$

$$\text{Need } f(\mathbf{x}) : \mathbb{R}^C \longrightarrow [0, \infty)^C, \quad \sum_{i=1}^C f(\mathbf{x})_c = 1$$

## Categorical distribution

$$p(y = c) = q_c, \quad y \in \{1 \dots C\}$$

$$\mathbf{q} \in \mathbb{R}^C \quad q_c \geq 0 \quad \forall c \in \{1 \dots C\} \quad \sum_{c=1}^C q_c = 1$$

$$p(y) = \prod q_c^{\mathbb{I}(y=c)} = q_y$$

$$\log p(y) =$$

## Softmax function

$$\text{softmax}(\mathbf{x})_c = \frac{e^{x_c}}{\sum_{j=1}^C e^{x_j}}$$

$$\text{softmax}(\mathbf{x}) = \begin{bmatrix} \frac{e^{x_1}}{\sum_{j=1}^C e^{x_j}} \\ \frac{e^{x_2}}{\sum_{j=1}^C e^{x_j}} \\ \vdots \\ \frac{e^{x_C}}{\sum_{j=1}^C e^{x_j}} \end{bmatrix}$$

## Softmax function

$$\text{softmax}(\mathbf{x})_c = \frac{e^{x_c}}{\sum_{j=1}^C e^{x_j}}$$

$$\sum_{i=1}^C \frac{e^{x_i}}{\sum_{j=1}^C e^{x_j}} = \frac{\sum_{i=1}^C e^{x_i}}{\sum_{j=1}^C e^{x_j}} = 1$$

$$\underset{c \in \{1, \dots, C\}}{\text{argmax } \mathbf{x}_c} = \underset{c \in \{1, \dots, C\}}{\text{argmax } \text{softmax}(\mathbf{x})_c}$$

## A probabilistic model for multi-class classification: Multinomial Logistic regression

$$y_i \sim \text{Categorical}(\text{softmax}(\mathbf{x}^T \mathbf{W}))$$

$$p(y_i = c \mid \mathbf{x}, \mathbf{W}) = \text{softmax}(\mathbf{x}^T \mathbf{W})_c = \frac{e^{\mathbf{x}^T \mathbf{w}_c}}{\sum_{j=1}^C e^{\mathbf{x}^T \mathbf{w}_j}}$$

# Maximum likelihood estimation for Multinomial Logistic regression

Model

$$y_i \sim \text{Categorical}(\text{softmax}(\mathbf{x}^T \mathbf{W}))$$

$$p(y_i = c \mid \mathbf{x}, \mathbf{W}) = \text{softmax}(\mathbf{x}^T \mathbf{W})_c = \frac{e^{\mathbf{x}^T \mathbf{w}_c}}{\sum_{j=1}^C e^{\mathbf{x}^T \mathbf{w}_j}}$$

Loss

$$\text{Loss}(\mathbf{W}) = \text{NLL}(\mathbf{W}, \mathbf{X}, \mathbf{y}) = - \sum_{i=1}^N \log p(y_i \mid \mathbf{x}_i, \mathbf{W})$$

$$\text{NLL}(\mathbf{W}, \mathbf{X}, \mathbf{y}) =$$

# Maximum likelihood estimation for Multinomial Logistic regression

Model

$$y_i \sim \text{Categorical}(\text{softmax}(\mathbf{x}^T \mathbf{W}))$$

$$p(y_i = c \mid \mathbf{x}, \mathbf{W}) = \text{softmax}(\mathbf{x}^T \mathbf{W})_c = \frac{e^{\mathbf{x}^T \mathbf{w}_c}}{\sum_{j=1}^C e^{\mathbf{x}^T \mathbf{w}_j}}$$

Loss

$$\text{Loss}(\mathbf{W}) = \text{NLL}(\mathbf{W}, \mathbf{X}, \mathbf{y}) = - \sum_{i=1}^N \log p(y_i \mid \mathbf{x}_i, \mathbf{W})$$

$$\text{NLL}(\mathbf{W}, \mathbf{X}, \mathbf{y}) =$$

$$\text{NLL}(\mathbf{W}, \mathbf{X}, \mathbf{y}) = -\sum_{i=1}^N \log \text{softmax}(\mathbf{x}_i^T \mathbf{W}^T)_{y_i} = -\sum_{i=1}^N \log \frac{e^{\mathbf{x}_i^T \mathbf{w}_{y_i}}}{\sum_{j=1}^C e^{\mathbf{x}_i^T \mathbf{w}_j}}$$

$$\text{NLL}(\mathbf{W}, \mathbf{X}, \mathbf{y}) = -\sum_{i=1}^N \left( \mathbf{x}_i^T \mathbf{w}_{y_i} - \log \sum_{j=1}^C e^{\mathbf{x}_i^T \mathbf{w}_j} \right)$$

-