# Generalized Linear Models

**Prepared by: Joseph Bakarji**

# **What if $y$ is a label?**
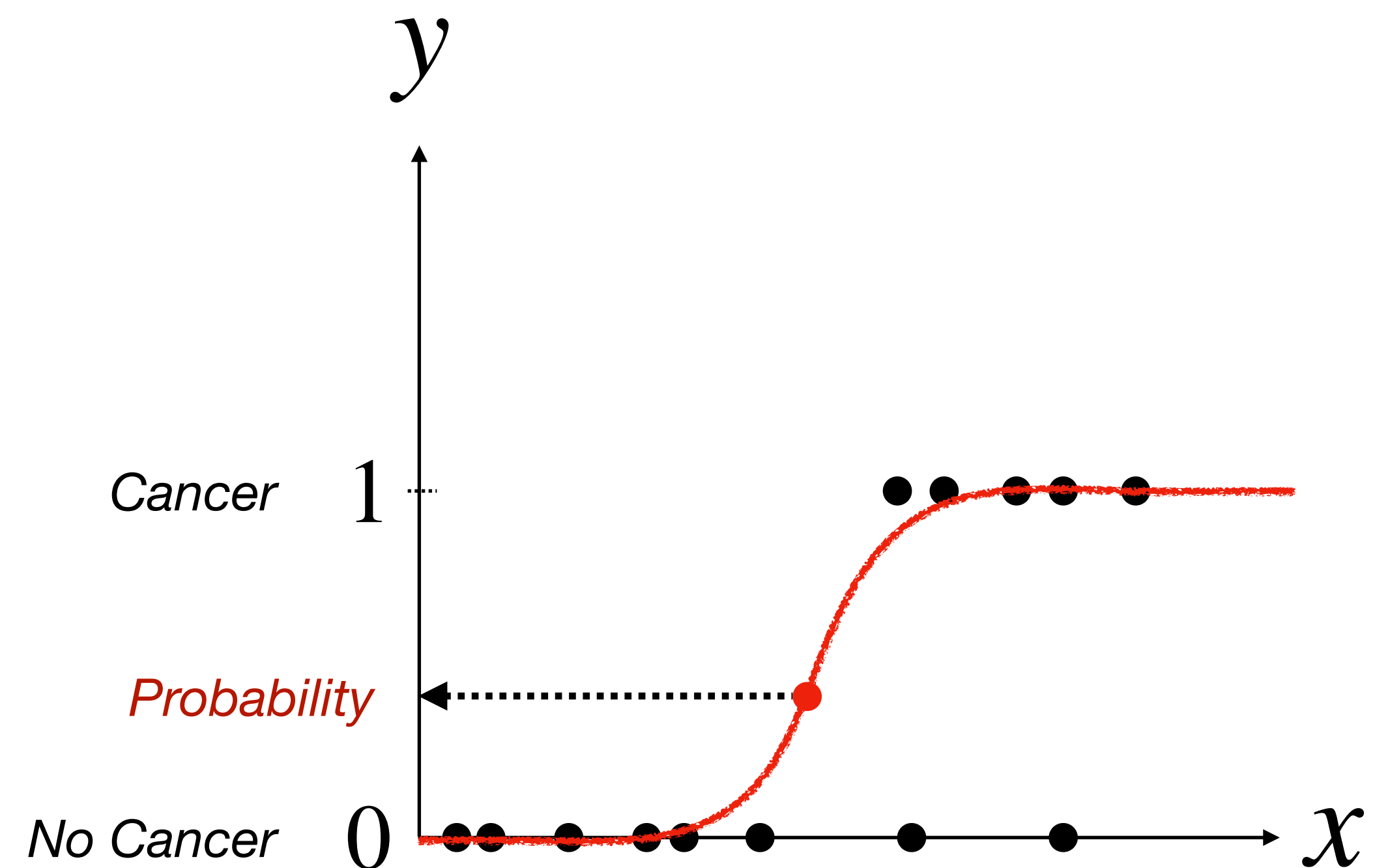
$$x \xrightarrow{\ h\ } y$$

Given the data,
find a **function** $h$,
that predicts $y$, given $x$
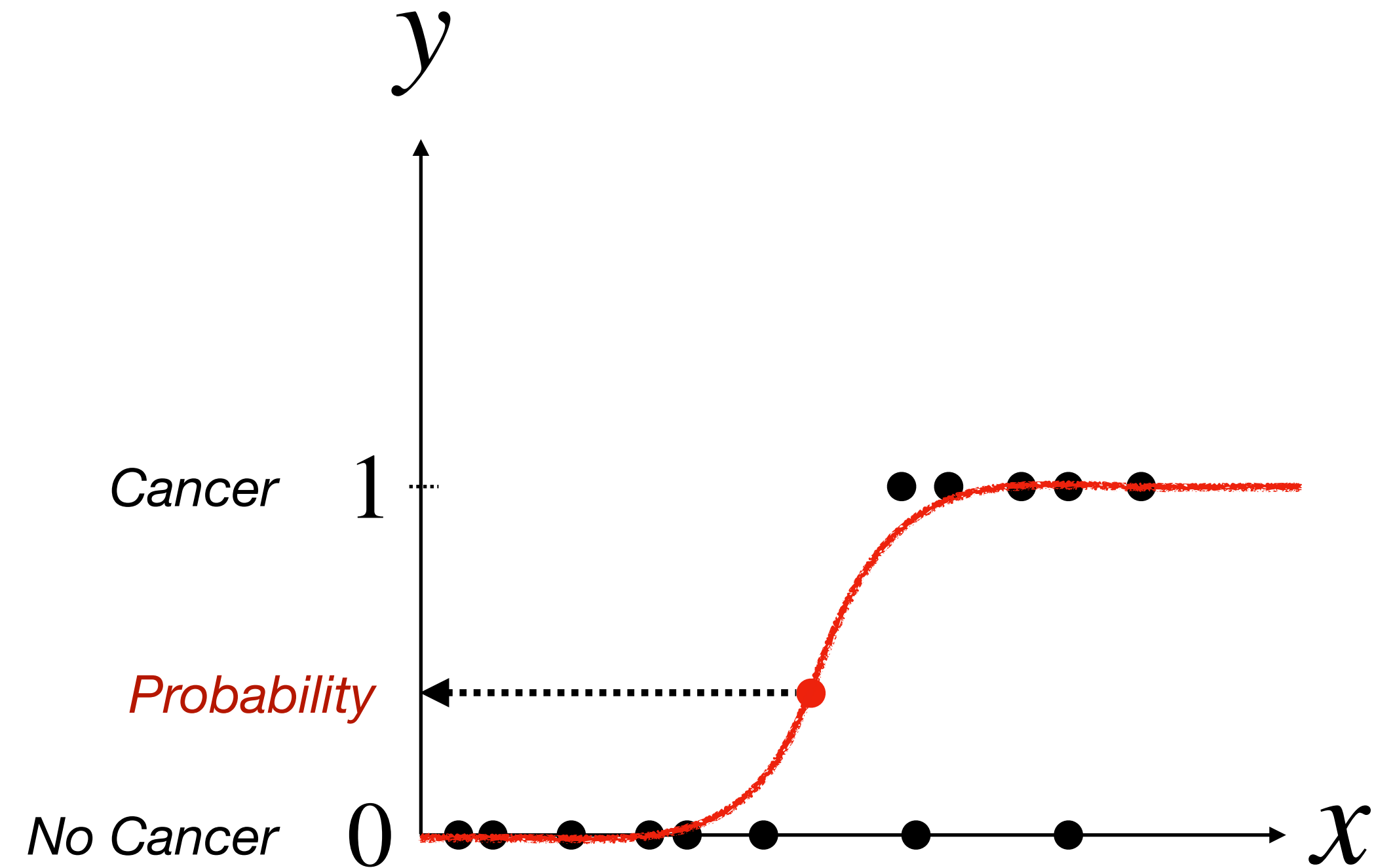
$$\mathbf{y} = h(\mathbf{x})$$

$$y \in [0,1]$$

**A smooth function that returns probability of occurrence**

# What if $y$ is a label?

$$y = h_\theta(x) \quad \& \quad y \in [0,1]$$

$$h_\theta(x) = \frac{1}{1 + e^{-(\theta^\top x)}}$$



1. **Define a predictor:** the logistic function ✅

2. **Define a loss:** distance between function and data **?**

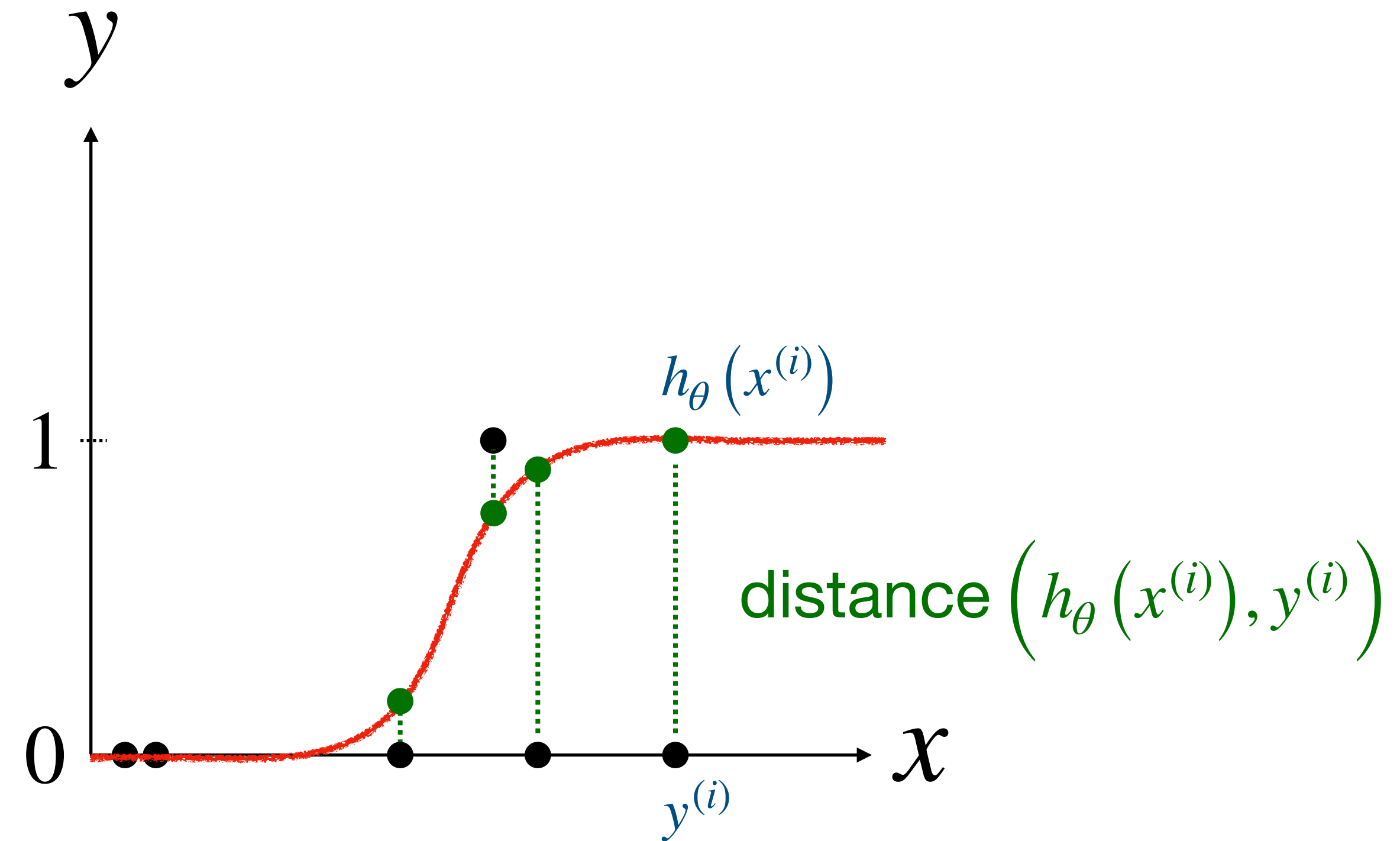3. **Optimize loss**

4. **Test model**

# Logistic Regression

$$y = h_\theta(x)$$

$$h_\theta(x) = \frac{1}{1 + e^{-(\theta^\top x)}} = g\left(\theta^\top x\right)$$



$h_\theta\left(x^{(i)}\right)$

$\text{distance}\left(h_\theta\left(x^{(i)}\right), y^{(i)}\right)$

$y^{(i)}$

Linear predictor
**negative log-likelihood or OLS**

$$J(\theta) = \frac{1}{2} \sum_{i=1}^{d} \left(h_\theta\left(x^{(i)}\right) - y^{(i)}\right)^2$$

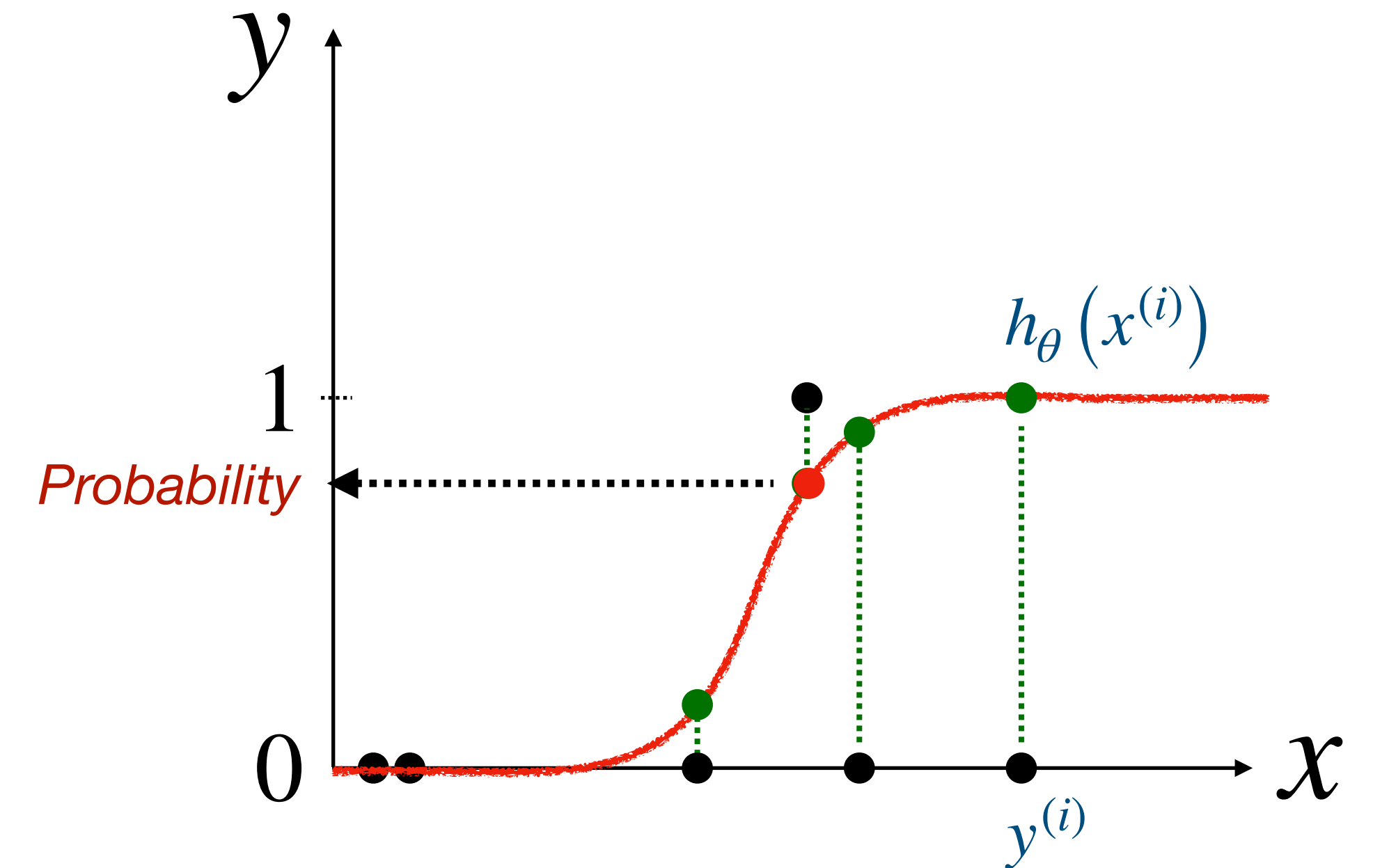Logistic predictor
**Binary-cross entropy loss**

$$\mathscr{L}(\theta) = \sum_{i=1}^{n} y^{(i)} \log h_\theta\left(x^{(i)}\right) + \left(1 - y^{(i)}\right) \log\left(1 - h_\theta\left(x^{(i)}\right)\right)$$

Gradient descent $\rightarrow$ Done!

# Why not Least Squares?

$$y = h_\theta(x)$$

$$h_\theta(x) = \frac{1}{1 + e^{-(\theta^\top x)}} = \sigma\left(\theta^\top x\right)$$



## Probability of output given input

$$P\left(y = 1 \,\middle|\, x; \theta\right) = h_\theta(x)$$

$$P\left(y = 0 \,\middle|\, x; \theta\right) = 1 - h_\theta(x)$$

$\longrightarrow$

True label

$$p(y \,|\, x; \theta) = \left(h_\theta(x)\right)^y \left(1 - h_\theta(x)\right)^{1-y}$$

**Likelihood!**

# Why not Least Squares?

$$y = h_\theta(x)$$

$$h_\theta(x) = \frac{1}{1 + e^{-(\theta^\top x)}} = \sigma\left(\theta^\top x\right)$$



$h_\theta\left(x^{(i)}\right)$

$h_\theta(x)$ *Probability*

$y^{(i)}$

## Probability of output given input

$$P\left(y = 1 \,\middle|\, x; \theta\right) = \sigma(\theta^\top x)$$

$$P\left(y = 0 \,\middle|\, x; \theta\right) = 1 - \sigma(\theta^\top x)$$

$\longrightarrow$

True label

$$p(y \,|\, x; \theta) = \left(\sigma(\theta^\top x)\right)^{y} \left(1 - \sigma(\theta^\top x)\right)^{1-y}$$

**Likelihood!**

# Maximize Log-likelihood

$$\mathcal{L}(\theta) = \sum_{i=1}^{n} y^{(i)} \log h_\theta \left( x^{(i)} \right) + \left( 1 - y^{(i)} \right) \log \left( 1 - h_\theta \left( x^{(i)} \right) \right)$$

## Update rule

**while** not converged:

$$\theta := \theta + \alpha \nabla_\theta \mathcal{L}(\theta)$$

**Derive** →

## Gradient Descent

**for** t = 1...T:

$$\theta := \theta - \alpha \sum_{i=1}^{n} \left( h_\theta \left( x^{(i)} \right) - y^{(i)} \right) x^{(i)}$$

🤯 **Same as linear regression** 🤔

# Generalized Linear Models

**Gaussian Distribution** $\longrightarrow$ **Linear Regression**

**Bernoulli Distribution** $\longrightarrow$ **Logistic Regression**

**Update rule**

$$\theta := \theta - \alpha \sum_{i=1}^{n} \left( h_\theta \left( x^{(i)} \right) - y^{(i)} \right) x^{(i)}$$

# Exponential Family
**Family of distributions for which we can derive <span style="color:darkred">the same update rule</span>**

<span style="color:darkred">**Assumption:**</span> $p(y \,|\, x; \theta)$ is an exponential family

**Data**

$$p(y; \eta) = b(y) \exp\left\{ \eta^\top y - a(\eta) \right\}$$

**Parameters**

- $b(y)$ is called the base measure (not depend on $\eta$)

- $a(\eta)$ is called the log partition function (not depend on $y$)

- $a(\eta)$, $y$ and $b(y)$ are scalar. $\eta$ and $y$ have the same dimensions.

# Example 1: Bernoulli Distribution -> Logistic Regression

**Data**

$$p(y; \eta) = b(y) \exp \left\{ \eta^\top y - a(\eta) \right\}$$

**Natural Parameters**

**Bernoulli Distribution**

$$p(y; \phi) = \phi^y (1 - \phi)^{1-y} = \exp \left\{ y \log \frac{\phi}{1 - \phi} + \log(1 - \phi) \right\}$$

$\eta$

$a(\eta)$

Show that term
is only a function of $\eta$

# Example 2: Gaussian Distribution -> Linear Regression

Data

$$p(y; \eta) = b(y) \exp \left\{ \eta^\top y - a(\eta) \right\}$$

Natural Parameters

## Gaussian Distribution

$$p(y; \mu) = \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{1}{2}(y - \mu)^2 \right\} = \underbrace{\frac{1}{\sqrt{2\pi}} e^{-y^2/2}}_{b(y)} \left\{ \underbrace{\mu}_{\eta} y - \underbrace{\frac{1}{2}\mu^2}_{a(\eta)} \right\}$$

# Why do we care?

**Data** $\leftarrow$ $\theta^\top x$

$$p(y; \eta) = b(y) \exp\left\{\eta^\top y - a(\eta)\right\}$$

**Natural** Parameters

**Inference is Easy:** $\quad E[y; \eta] = \dfrac{da(\eta)}{d\eta} \qquad\qquad Var[y; \eta] = \dfrac{d^2 a(\eta)}{d\eta^2}$

**Learning is Easy:** $\quad$ **Maximum Likelihood Estimation** leads to **convex** problem in $\eta$

# Generalized Linear Models

**Assumption:** $p(y \mid x; \theta)$ is an exponential family

**Data Type $\rightarrow$ Probability Distribution**

Binary $\rightarrow$ Bernoulli $\longrightarrow$ **Logistic Regression**

Real $\rightarrow$ Gaussian $\longrightarrow$ **Linear Regression**

Counts $\rightarrow$ Poisson

Positive Real $\rightarrow$ Gamma, Exponential

Distributions $\rightarrow$ Dirichlet

# Generalized Linear Models

**Assumption:** $p(y \mid x; \theta)$ is an exponential family

The natural parameter is linear in the inputs

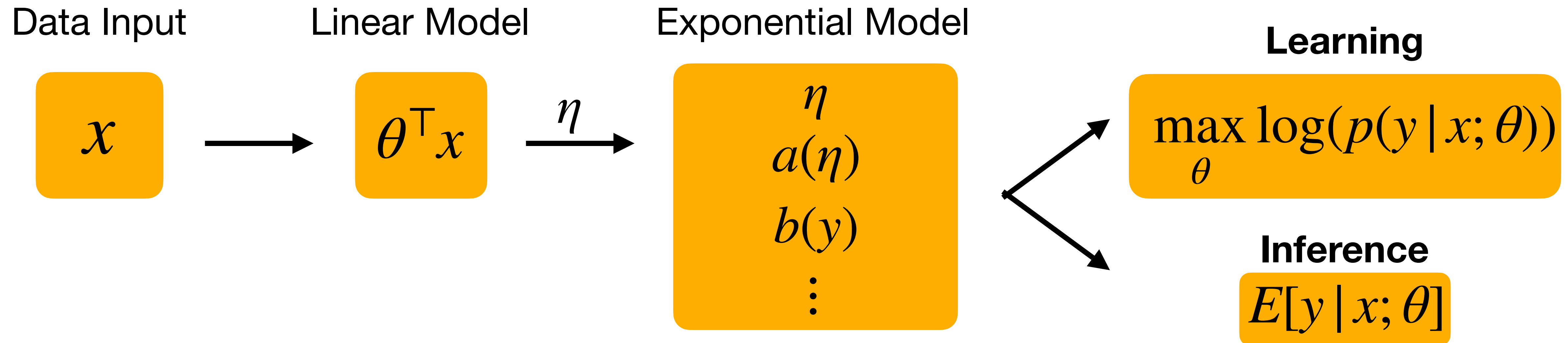$$\eta = \theta^\top x$$

Predictor is a natural consequence

$$h_\theta(x) = E[y \mid x; \theta]$$
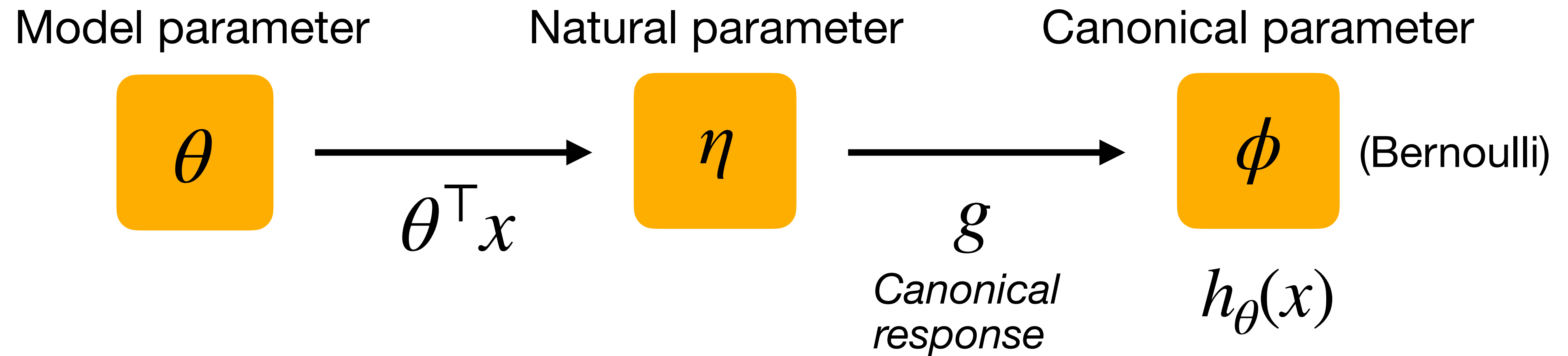
# Generalized Linear Models

**Assumption:** $p(y \mid x; \theta)$ is an exponential family

Data Input

Linear Model

Exponential Model

$$x \longrightarrow \theta^\top x \xrightarrow{\eta}$$

$$\begin{array}{c} \eta \\ a(\eta) \\ b(y) \\ \vdots \end{array}$$

**Learning**

$$\max_{\theta} \log(p(y \mid x; \theta))$$

**Inference**

$$E[y \mid x; \theta]$$

**Update Rule:**

$$\theta := \theta - \alpha \sum_{i=1}^{n} \left( h_\theta \left( x^{(i)} \right) - y^{(i)} \right) x^{(i)}$$
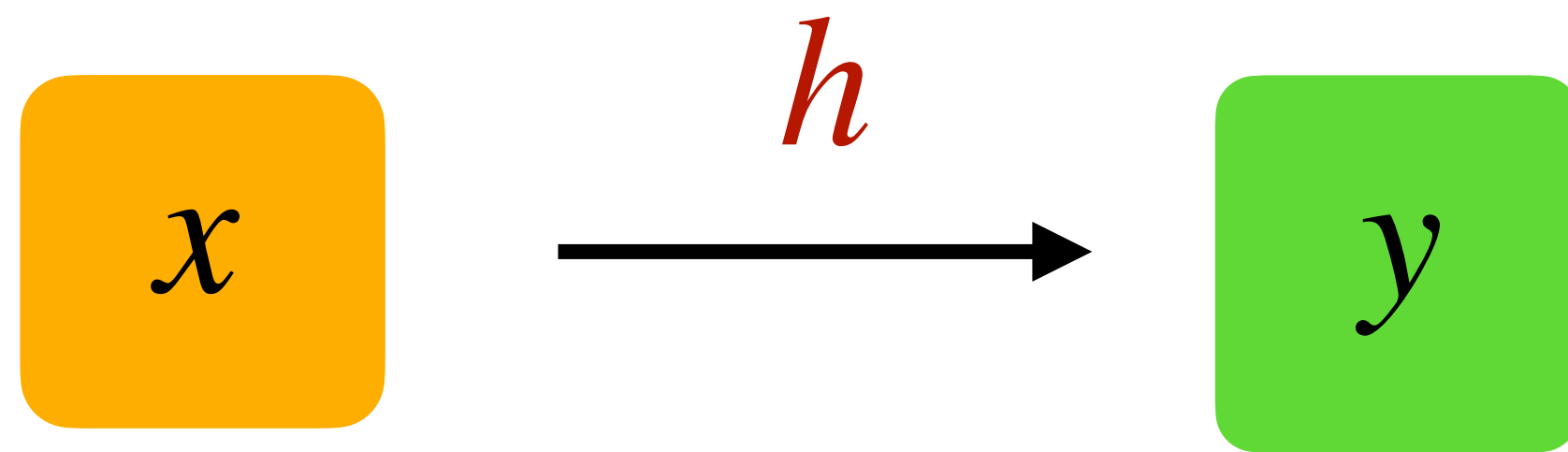
# Terminology

| Model parameter | | Natural parameter | | Canonical parameter |
|:---:|:---:|:---:|:---:|:---:|
| $\theta$ | $\xrightarrow{\quad\theta^{\top}x\quad}$ | $\eta$ | $\xrightarrow[\substack{\textit{Canonical} \\ \textit{response}}]{\quad g\quad}$ | $\phi$ (Bernoulli) |
| | | | | $h_\theta(x)$ |

**Logistic Regression:** $\qquad h_\theta(x) = E[y \mid x; \theta] \qquad\qquad \phi = \dfrac{1}{1 + e^{-\eta}} = \dfrac{1}{1 + e^{-\theta^{\top}x}}$

# Back to classification

**What if we have more outputs?**

$x$   $h$   $y$

$y = h_\theta(x)$

$y \in [0,1]$

$y$

Cancer   1

Probability

No Cancer   0

# Classification

$$x = \left[ x_1, x_2 \right]$$

| $x_1$ | $x_2$ | $y$ |
|-------|-------|-----|
| -2    | -1    | ● |
| 3     | 1     | ★ |
| 2     | 3     | ● |
| 1     | -1    | ★ |
| ⋮     |       |   |

★ 1

● 0

$x_2$

*Decision boundary*

$x_1$

$\theta$

$\theta^\top x = 0$

*Logistic Regression*

$$h_\theta(x) = \sigma(\theta^\top x)$$

*how confident?*

**score**

$\theta^\top x$

*how correct?*

**margin**

$(\theta^\top x)y$

For $y \in [1, -1]$

# **Multiclass classification** - Softmax

$$\mathbf{y} = h_\theta(\mathbf{x})$$



$k$ discrete values for representing output



**One-hot encoding**

$$\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} \qquad \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \qquad \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix} \qquad \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$
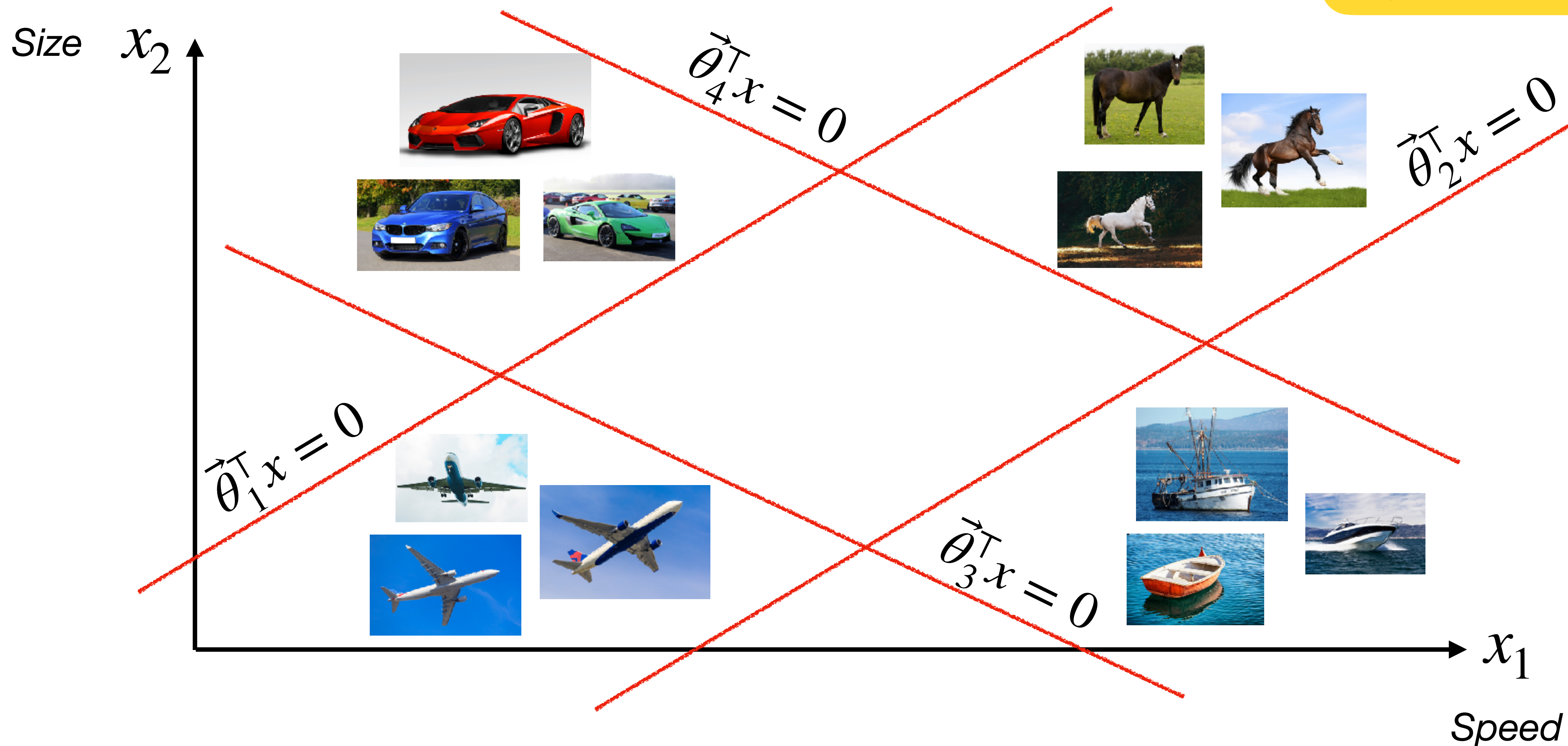
car        plane        boat        horse

# Multi-class classification - Softmax
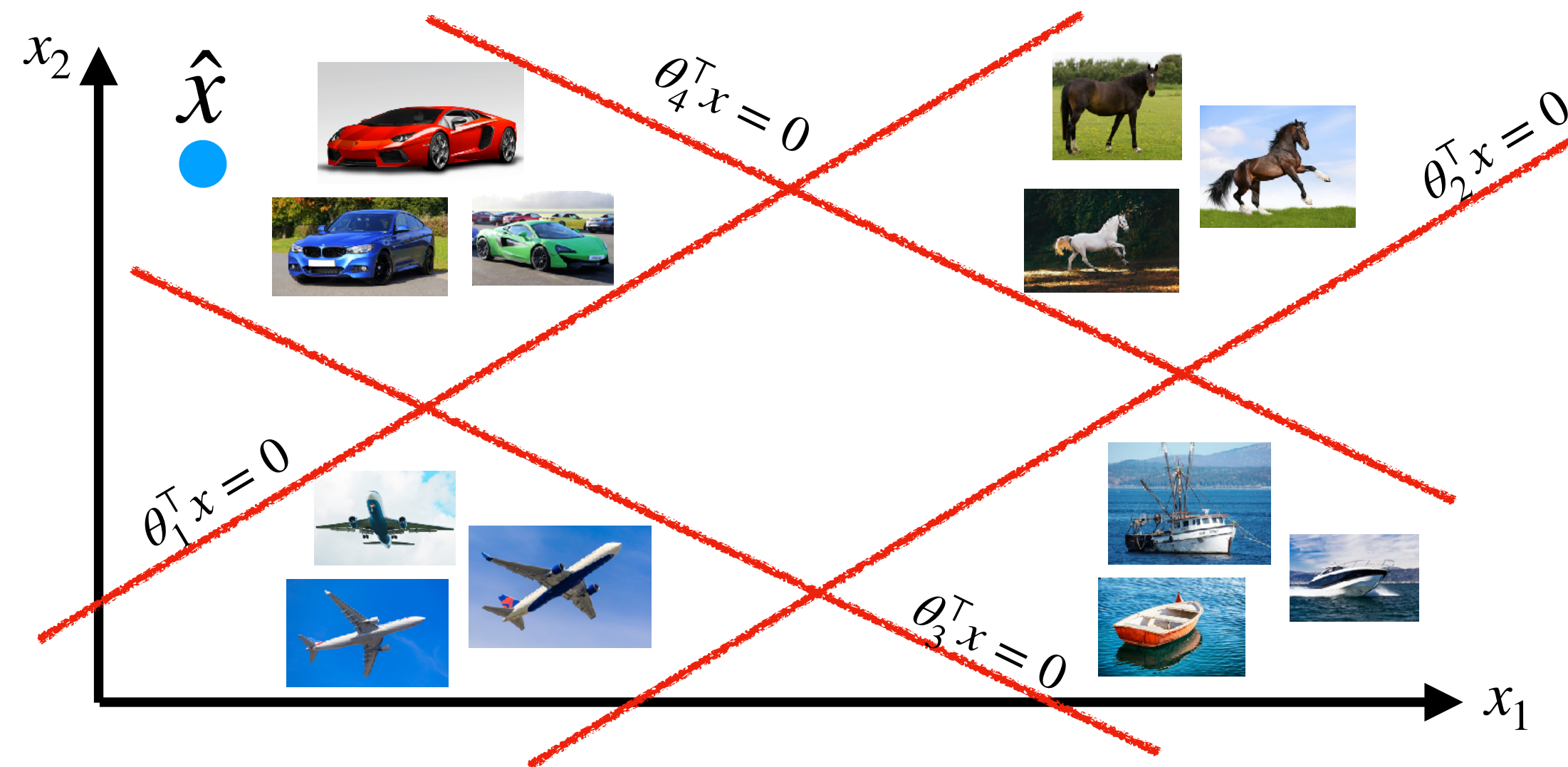
**WARNING!!!**
**Notation Alert**
$\vec{\theta}_i$ **is a vector**

*Size* $x_2$

$\vec{\theta}_4^\top x = 0$

$\vec{\theta}_2^\top x = 0$

$\vec{\theta}_1^\top x = 0$

$\vec{\theta}_3^\top x = 0$

$x_1$

*Speed*

# How to turn scores into probabilities?



WARNING!!!

Notation Alert

$\vec{\theta}_i$ is a vector

**Score**

$$\vec{\theta}_1^\top \hat{x} = 3$$

$$\vec{\theta}_2^\top \hat{x} = -0.3$$

$$\vec{\theta}_3^\top \hat{x} = -0.8$$

$$\vec{\theta}_4^\top \hat{x} = -22$$

$\xrightarrow{\exp}$

**Positive Measure**

$$\exp(3) = 20.1$$

$$\exp(-0.3) = 0.75$$

$$\exp(-0.8) = 0.2$$

$$\exp(-22) = 0.00..1$$

**Normalize** $\longrightarrow$

**Softmax**

$$\hat{p}(y = i \mid x; \theta) = \frac{\exp\left(\vec{\theta}_i^\top x\right)}{\sum_{j=1}^{k} \exp\left(\vec{\theta}_j^\top x\right)}$$

# How do you train?

**Given Label**

$$\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$x$

car

$p$

1

Probability

Label smoothing

car  plane  boat  horse

**Inference (prediction)**

$\hat{p}$

1

Probability

car  plane  boat  horse

$$\textbf{min} \ \ \text{CrossEntropy}(p, \hat{p}) = -\sum_{i=1}^{k} p(y = i) \log\left(\hat{p}(y = i)\right)$$

$$= -\log\left(\hat{p}(y = 1)\right)$$

# How do you train?



**Inference (prediction)**

$\hat{p}$

Probability

1

car    plane    boat    horse

1      2      3      4

$$\text{CrossEntropy}(p, \hat{p}) = -\sum_{i=1}^{k} p(y = i) \log\left(\hat{p}(y = i)\right)$$

**Ground Truth**

$$Logit \quad = -\log\left(\hat{p}(y = 1)\right)$$

$$= -\log\left(\frac{\exp\left(\vec{\theta}_i^\top x\right)}{\sum_{j=1}^{k} \exp\left(\vec{\theta}_j^\top x\right)}\right)$$

**Train with Gradient Descent!**