

Maximum Likelihood Estimation & Logistic Regression

Slides adapted from David Sontag and Andrew Ng.

Maximum Likelihood Estimation (MLE)

Framework:

- Observed data D (observations)
- Hypothesize data has a specific probability distribution parameterized by unknown parameter values θ : i.e., distribution $P_{\theta}(D)$ is known
- Goal: estimate (learn) the parameter values θ .
- MLE: Choose parameter values θ that maximize $P_{\theta}(D)$

Thumbtack example

- $P_{\theta}(\text{Heads}) = \theta$, $P_{\theta}(\text{Tails}) = 1 - \theta$. We want to estimate θ from observational data.



- Make observations:

$$D = \{y_i \mid i=1, \dots, m\} \quad y_i = \text{H or T}$$

- Need a model $P_{\theta}(D)$
- Let α_H be number of heads in D ; α_T be number of tails in D .
- Natural model (Why?)

$$P_{\theta}(D) = \theta^{\alpha_H} (1 - \theta)^{\alpha_T}$$

Called the “likelihood” of the data under the model.

Maximum Likelihood Estimation

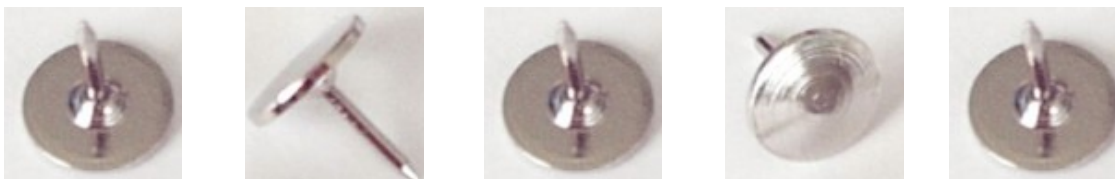
- **Data:** Observed set D : sequence of heads and tails, with α_H Heads and α_T Tails.
- **Model:** $P_{\theta}(D) = \theta^{\alpha_H} (1 - \theta)^{\alpha_T}$
- **Learning:** find θ that maximizes the probability of the observation D , i.e., find:

$$\hat{\theta} = \arg \max_{\theta} P_{\theta}(D)$$

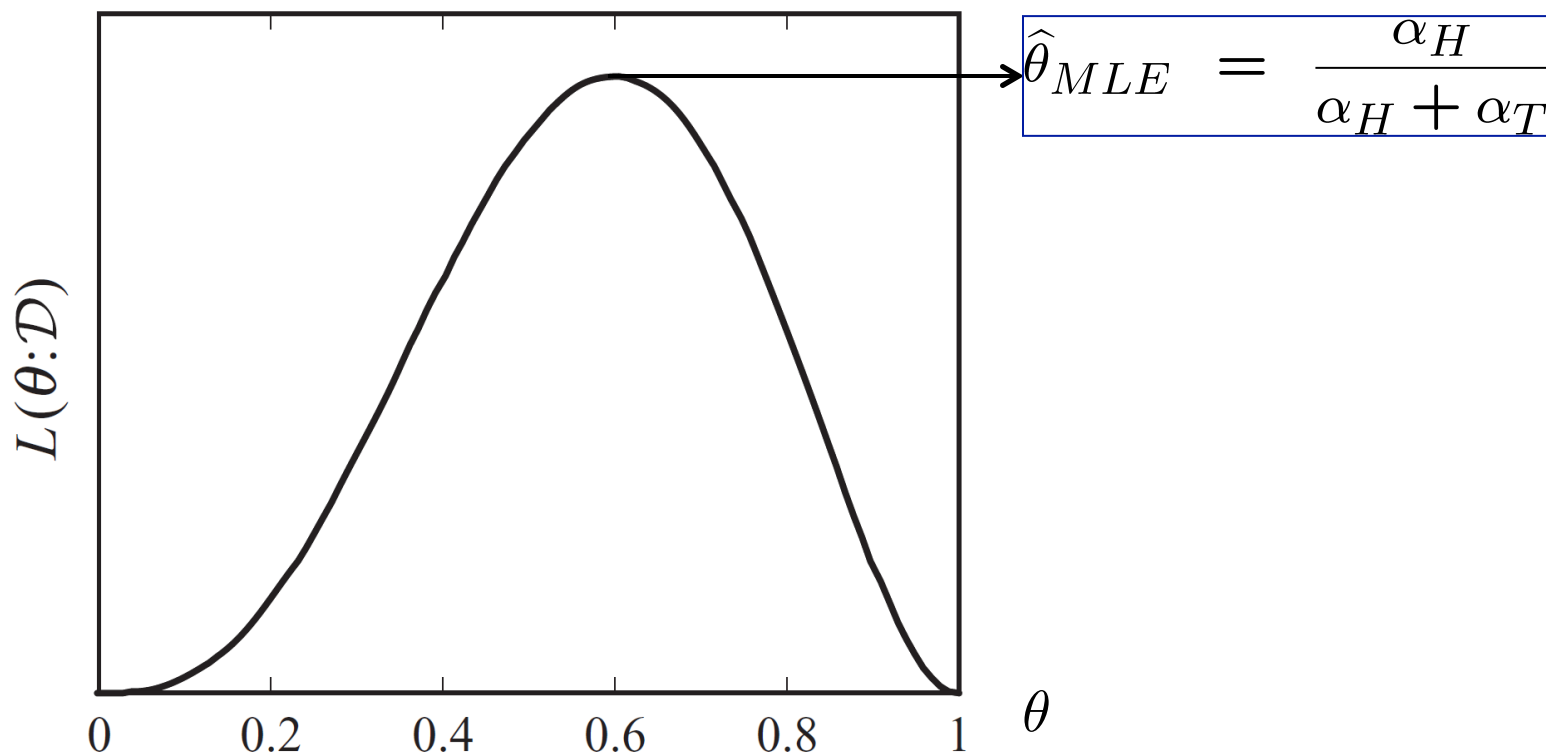
- **Taking derivative and setting to zero, get:**

$$\boxed{\hat{\theta}_{MLE} = \frac{\alpha_H}{\alpha_H + \alpha_T}}$$

Data



$$L(\boldsymbol{\theta}; \mathcal{D}) = P_{\boldsymbol{\theta}}(\mathcal{D})$$



Logistic Regression

- Popular type of supervised machine learning for classification
- Classification, not regression!
- Gives probabilities for classification, e.g., email is spam with probability 0.86
- Can be viewed as a MLE estimator
- Often used in neural networks

Example:

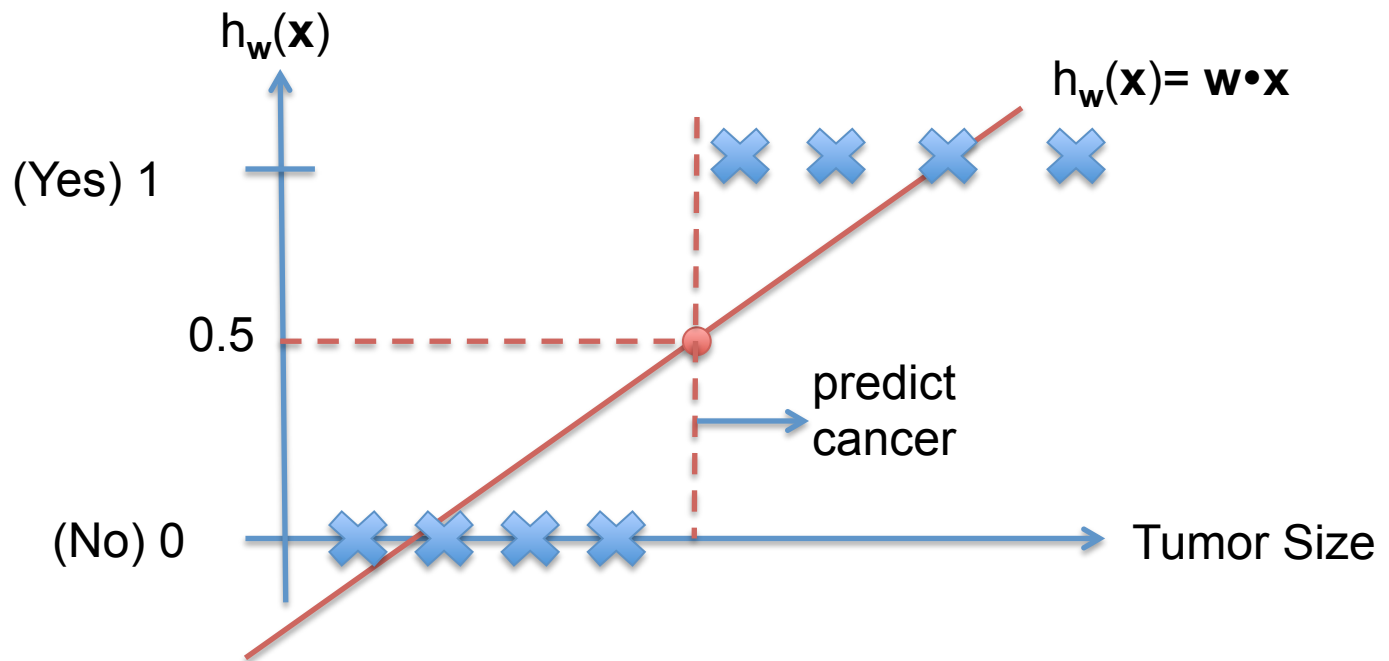
- Tumor: Malignant/ Benign?

$y \text{ in } \{0,1\}$

0: “Negative Class” (e.g., benign tumor)

1: “Positive Class” (e.g., malignant tumor)

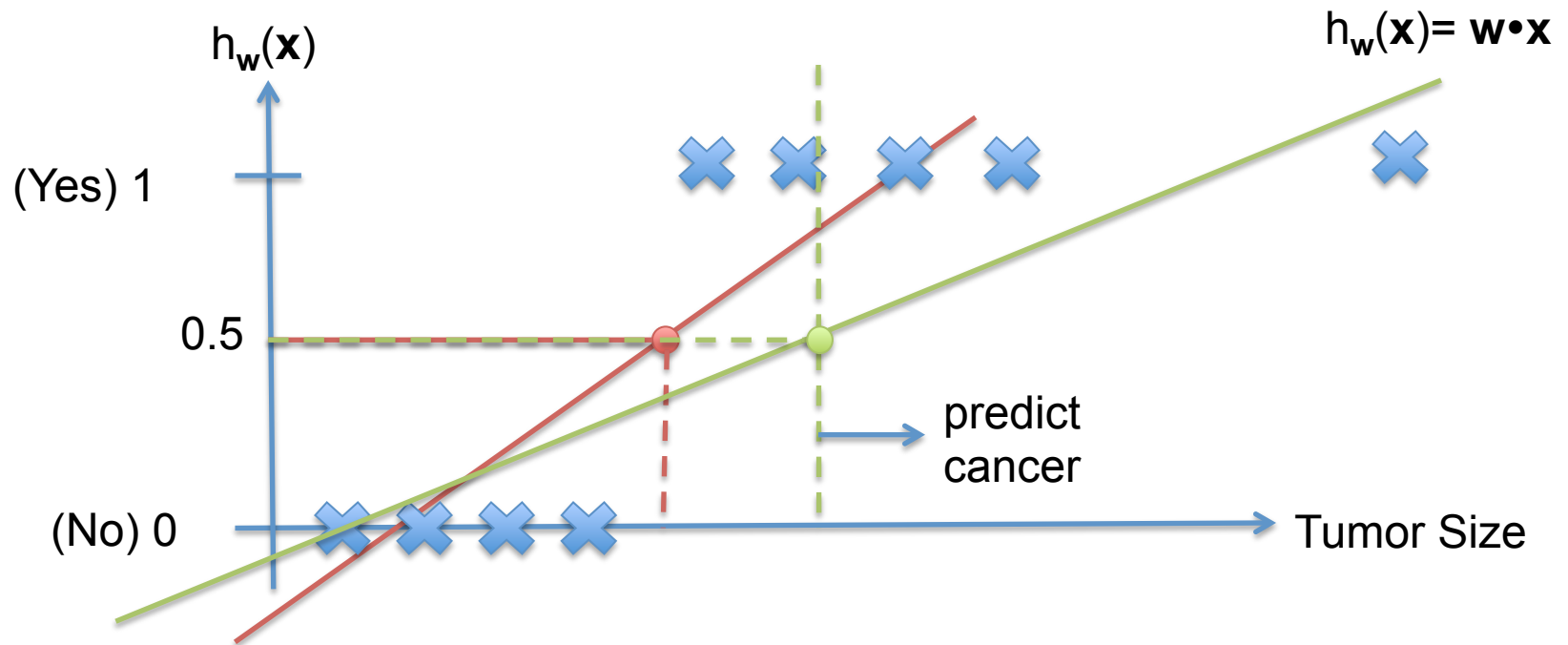
Let's try to predict with ordinary regression



Natural threshold classifier:

- If $h_w(\mathbf{x}) \geq 0.5$, predict “y=1”
- If $h_w(\mathbf{x}) < 0.5$, predict “y=0”

Additional data point



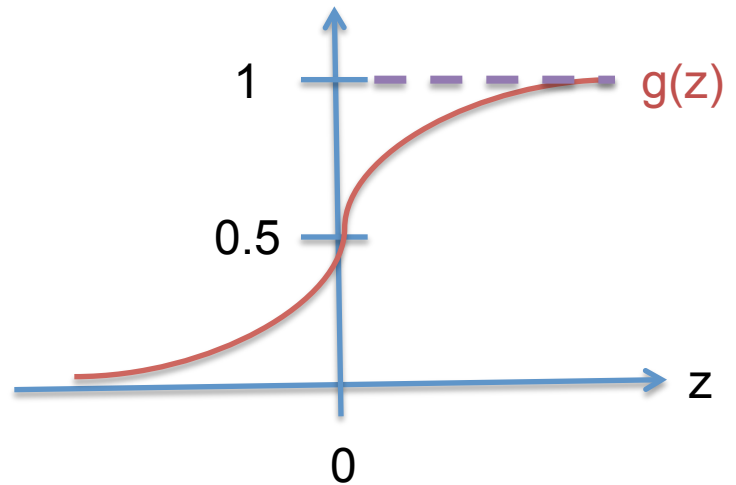
Linear regression with natural 0.5 threshold does not look good here.

Graphically, what kind of function would be a good fit?

Sigmoid function

$$g(z) = \frac{1}{1 + e^{-z}}$$

Sigmoid function = Logistic function

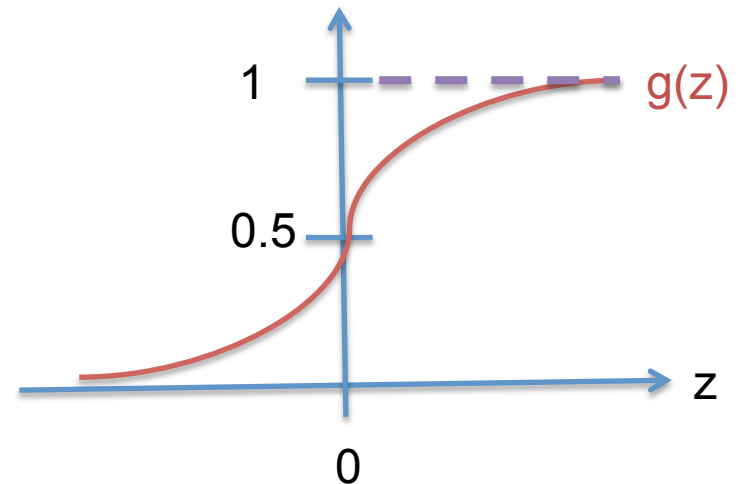


Logistic Regression Model

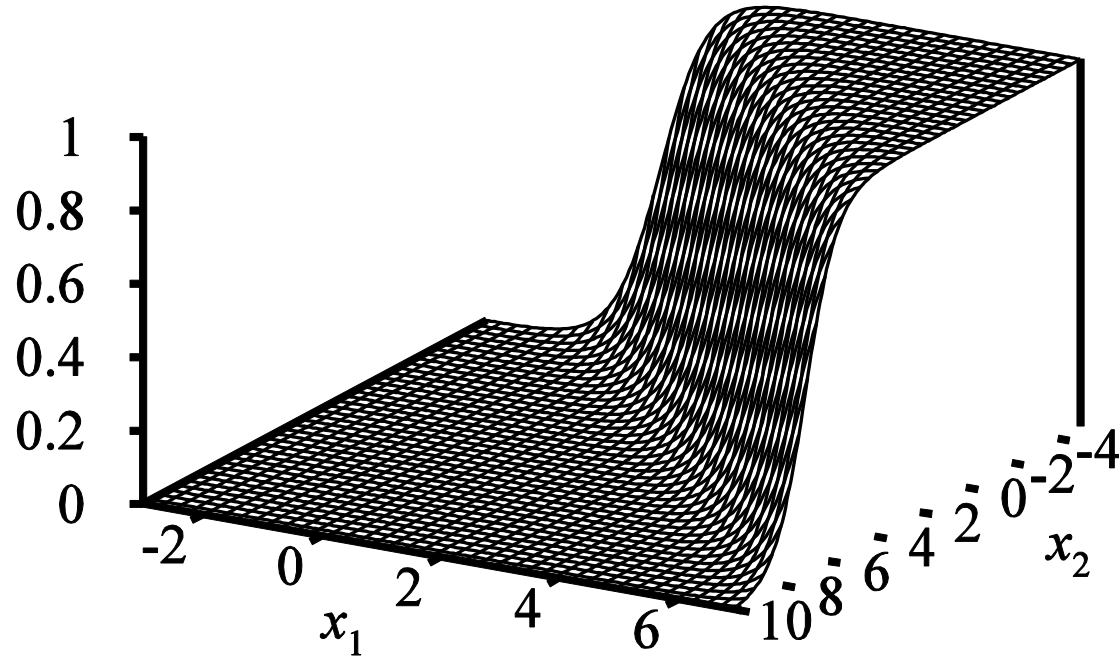
$$g(z) = \frac{1}{1 + e^{-z}}$$

$$h_{\mathbf{w}}(\mathbf{x}) = g(\mathbf{w} \cdot \mathbf{x})$$

- Note that $0 \leq h_{\mathbf{w}}(\mathbf{x}) \leq 1$
- Predict $y = 1$ if $h_{\mathbf{w}}(\mathbf{x}) > \frac{1}{2}$; otherwise predict $y = 0$.
- $h_{\mathbf{w}}(\mathbf{x})$ can be interpreted as a probability, eg., probability of cancer
- Can choose \mathbf{w} to optimize the fit to data (later).



Logistic Function in n Dimensions



Tumor example

- Suppose we have learned \mathbf{w} . Observe \mathbf{x} for new patient and want to predict if patient has cancer
- $h_{\mathbf{w}}(\mathbf{x}) = g(\mathbf{w} \bullet \mathbf{x})$ estimated probability that patient has cancer
- Example: $\mathbf{x} = (x_0, x_1) = (1, \text{tumorSize})$
- Suppose $h_{\mathbf{w}}(\mathbf{x}) = 0.7$.
- Can tell patient 70% chance tumor is malignant
- That is, probability model $P_{\mathbf{w}}(y=1 | \mathbf{x}) = h_{\mathbf{w}}(\mathbf{x})$

Summary

- Use labeled data to learn \mathbf{w}
- Observe new \mathbf{x}
- Given \mathbf{x} , we say $y=1$ with estimated probability $h_{\mathbf{w}}(\mathbf{x}) = g(\mathbf{w} \bullet \mathbf{x})$, where

$$g(z) = \frac{1}{1 + e^{-z}}$$

- Alternatively way of saying it: $P_{\mathbf{w}}(y=1 | \mathbf{x}) = h_{\mathbf{w}}(\mathbf{x})$.
- But how do we learn \mathbf{w} ?

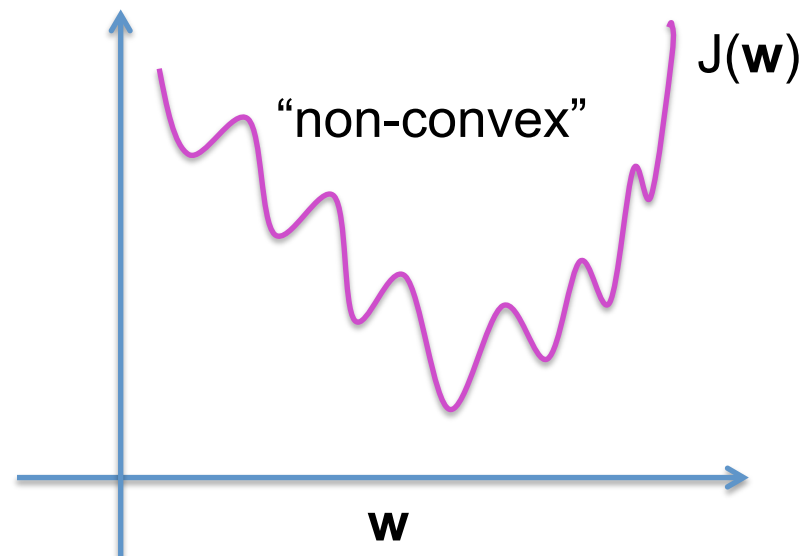
Learning the Parameters \mathbf{w}

- Training set: $\{(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \dots, (\mathbf{x}^{(m)}, y^{(m)})\}$
- How about choosing \mathbf{w} to minimize MSE (as usual)?

$$J(\mathbf{w}) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(h_{\mathbf{w}}(\mathbf{x}^{(i)}), y^{(i)})$$

$$\text{Cost}(h_{\mathbf{w}}(\mathbf{x}), y) = \frac{1}{2} (h_{\mathbf{w}}(\mathbf{x}) - y)^2$$

$$h_{\mathbf{w}}(x) = \frac{1}{1 + e^{-\mathbf{w} \cdot \mathbf{x}}}$$



Learning the Parameters \mathbf{w}

- Instead try using MLE. Find \mathbf{w} that maximizes the probability of the observation. Maximize:

$$P_{\mathbf{w}}(y^{(1)}, y^{(2)}, \dots, y^{(m)} | \mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(m)})$$

- Need model for $P_{\mathbf{w}}(y^{(1)}, y^{(2)}, \dots, y^{(m)} | \mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(m)})$:
 - Assume each observed data point is conditionally independent:

$$P_{\mathbf{w}}(y^{(1)}, y^{(2)}, \dots, y^{(m)} | \mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(m)}) = \\ P_{\mathbf{w}}(y^{(1)} | \mathbf{x}^{(1)}) \times P_{\mathbf{w}}(y^{(2)} | \mathbf{x}^{(2)}) \times \dots \times P_{\mathbf{w}}(y^{(m)} | \mathbf{x}^{(m)})$$

- Assume logistic function probabilities:

$$P_{\mathbf{w}}(y^{(i)}=1 | \mathbf{x}^{(i)}) = h_{\mathbf{w}}(\mathbf{x}^{(i)})$$

$$P_{\mathbf{w}}(y^{(i)}=0 | \mathbf{x}^{(i)}) = 1 - h_{\mathbf{w}}(\mathbf{x}^{(i)})$$

$$h_{\mathbf{w}}(\mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w} \cdot \mathbf{x}}}$$

Choose \mathbf{w} to maximize:

- $P_{\mathbf{w}}(y^{(1)}, y^{(2)}, \dots, y^{(m)} | \mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(m)})$
- Because log is an increasing function, we can instead maximize
$$\begin{aligned} \log (P_{\mathbf{w}}(y^{(1)}, y^{(2)}, \dots, y^{(m)} | \mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(m)})) &= \\ \log [P_{\mathbf{w}}(y^{(1)} | \mathbf{x}^{(1)}) \times P_{\mathbf{w}}(y^{(2)} | \mathbf{x}^{(2)}) \times \dots \times P_{\mathbf{w}}(y^{(m)} | \mathbf{x}^{(m)})] &= \\ \log P_{\mathbf{w}}(y^{(1)} | \mathbf{x}^{(1)}) + \log P_{\mathbf{w}}(y^{(2)} | \mathbf{x}^{(2)}) + \dots + \log P_{\mathbf{w}}(y^{(m)} | \mathbf{x}^{(m)}) \end{aligned}$$
- $\log P_{\mathbf{w}}(y^{(i)} = 1 | \mathbf{x}^{(i)}) = \log h_{\mathbf{w}}(\mathbf{x}^{(i)})$
 $\log P_{\mathbf{w}}(y^{(i)} = 0 | \mathbf{x}^{(i)}) = \log (1 - h_{\mathbf{w}}(\mathbf{x}^{(i)}))$
- So $P_{\mathbf{w}}(y^{(i)} | \mathbf{x}^{(i)}) = y^{(i)} \log h_{\mathbf{w}}(\mathbf{x}^{(i)}) + (1 - y^{(i)}) \log (1 - h_{\mathbf{w}}(\mathbf{x}^{(i)}))$

$$\begin{aligned} J(\mathbf{w}) &= \frac{1}{m} \sum_{i=1}^m \text{Cost}(h_{\mathbf{w}}(\mathbf{x}^{(i)}), y^{(i)}) \\ &= -\frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \log h_{\mathbf{w}}(\mathbf{x}^{(i)}) + (1 - y^{(i)}) \log (1 - h_{\mathbf{w}}(\mathbf{x}^{(i)})) \right] \end{aligned}$$

Convex
function!

Summary

- Want to find a \mathbf{w} so that the logistic regression fits the data.
- Usual MSE error cost function leads to non-convex optimization problem.
- Instead consider finding \mathbf{w} that maximizes likelihood (MLE). This is equivalent to finding \mathbf{w} that minimizes:

$$J(\mathbf{w}) = -\frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \log h_{\mathbf{w}}(\mathbf{x}^{(i)}) + (1 - y^{(i)}) \log (1 - h_{\mathbf{w}}(\mathbf{x}^{(i)})) \right]$$

- Convex optimization problem

Gradient Descent

$$J(\mathbf{w}) = -\frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \log h_{\mathbf{w}}(\mathbf{x}^{(i)}) + (1 - y^{(i)}) \log (1 - h_{\mathbf{w}}(\mathbf{x}^{(i)})) \right]$$

Want $\min_{\mathbf{w}} J(\mathbf{w})$:

Repeat {

$$w_j := w_j - \alpha \frac{\partial}{\partial w_j} J(\mathbf{w})$$

(simultaneously update all \mathbf{w}_j)

}

$$\frac{\partial}{\partial w_j} J(\mathbf{w}) = \frac{1}{m} \sum_{i=1}^m (h_{\mathbf{w}}(\mathbf{x}^{(i)}) - y^{(i)}) x_j^{(i)}$$

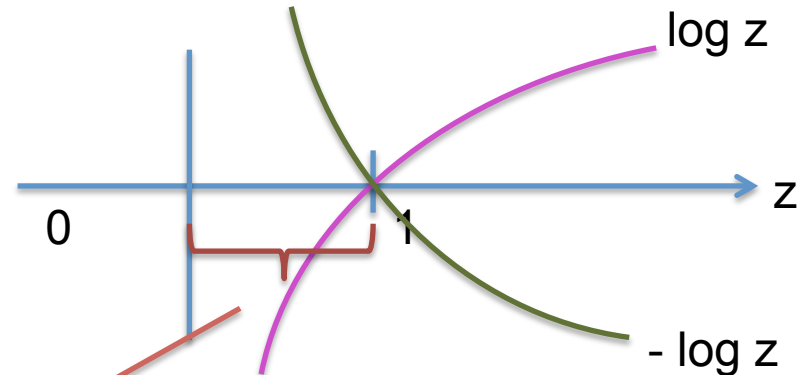
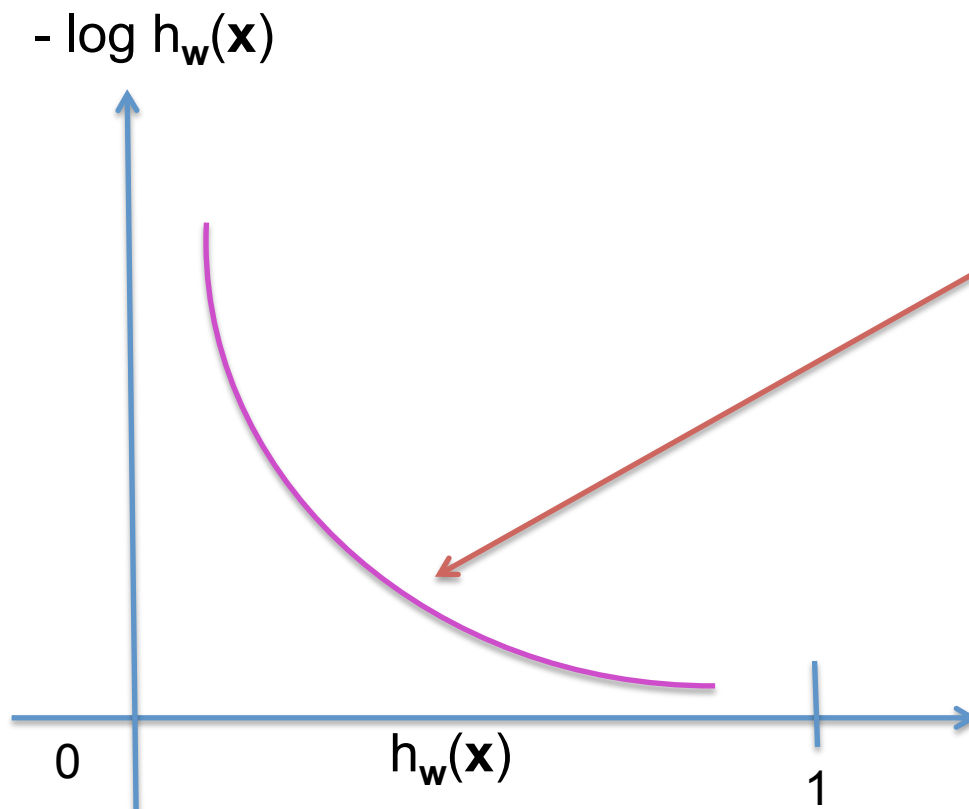
$$w_j := w_j - \alpha \sum_{i=1}^m (h_{\mathbf{w}}(\mathbf{x}^{(i)}) - y^{(i)}) x_j^{(i)}$$

Algorithm looks identical to linear regression! But it is not!

Some intuition into cost function

$$\text{Cost}(h_w(\mathbf{x}), y) = \begin{cases} -\log(h_w(\mathbf{x})) & \text{if } y = 1 \\ -\log(1 - h_w(\mathbf{x})) & \text{if } y = 0 \end{cases}$$

Consider $y=1$



As $h_w(\mathbf{x}) \rightarrow 1$, Cost $\rightarrow 0$
But as $h_w(\mathbf{x}) \rightarrow 0$, Cost $\rightarrow \infty$
Captures intuition that if $h_w(\mathbf{x}) = 0$, (predict $P_w(y=1 | \mathbf{x})=0$), we'll penalize learning algorithm by a very large cost.

Summary: Logistic regression cost function

$$\begin{aligned} J(\mathbf{w}) &= \frac{1}{m} \sum_{i=1}^m \text{Cost}(h_{\mathbf{w}}(\mathbf{x}^{(i)}), y^{(i)}) \\ &= -\frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \log h_{\mathbf{w}}(\mathbf{x}^{(i)}) + (1 - y^{(i)}) \log (1 - h_{\mathbf{w}}(\mathbf{x}^{(i)})) \right] \end{aligned}$$

To fit parameters \mathbf{w} :

$$\min_{\mathbf{w}} J(\mathbf{w})$$

To make a prediction given new \mathbf{x} :

$$\text{Output} \quad h_{\mathbf{w}}(x) = \frac{1}{1 + e^{-\mathbf{w} \cdot \mathbf{x}}}$$