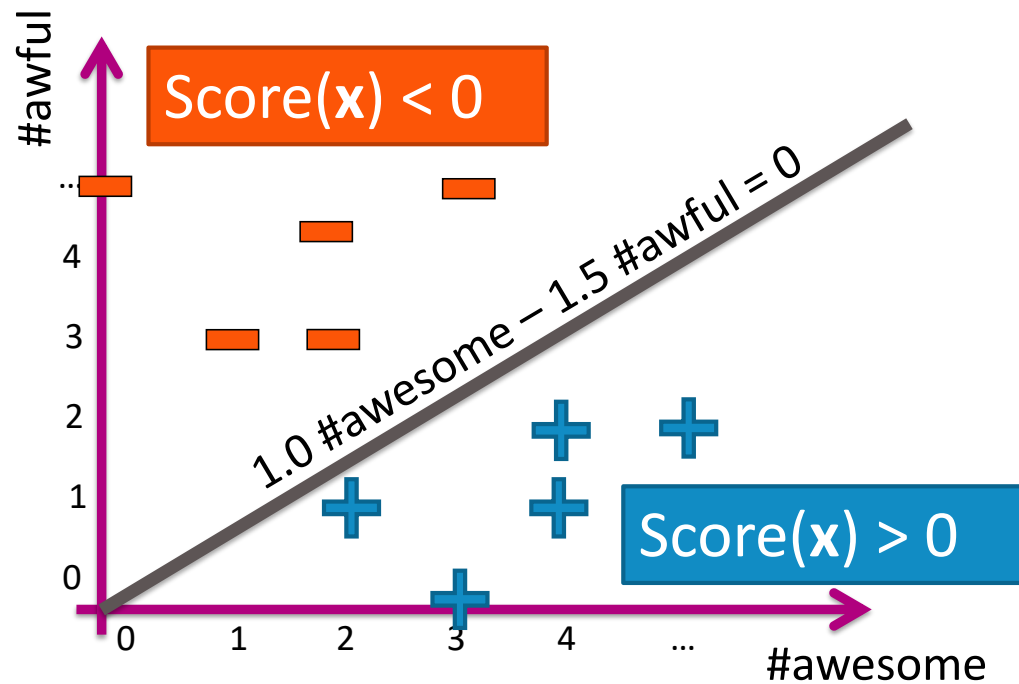


Review Classification

Linear classifier

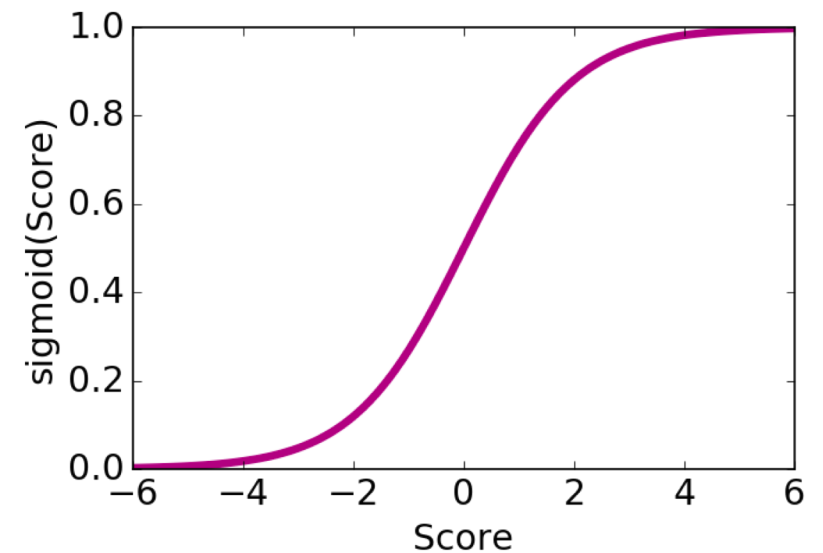
$$\begin{aligned}\text{Score}(x_i) &= w_0 h_0(x_i) + w_1 h_1(x_i) + \dots + w_D h_D(x_i) \\ &= \mathbf{w}^T \mathbf{h}(x_i)\end{aligned}$$



Logistic function (sigmoid, logit)

$$\text{sigmoid}(\text{Score}) = \frac{1}{1 + e^{-\text{Score}}}$$

Score	$-\infty$	-2	0.0	+2	$+\infty$
sigmoid(Score)					



Probabilistic Model

Distribution over sales prices

For houses with a given # sq.ft. () , what house prices \$ are we likely to see?

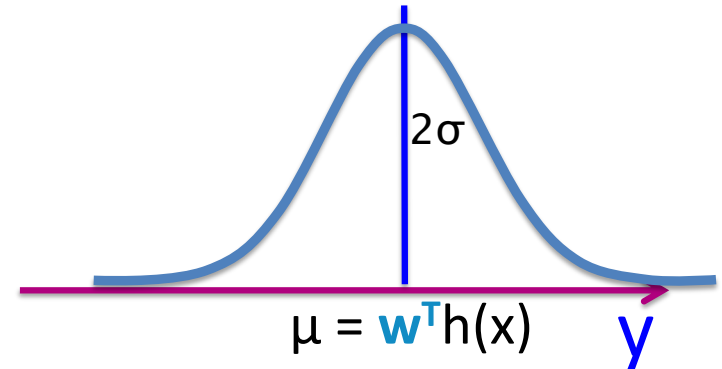


Compare and contrast regression models

- Linear regression with Gaussian errors

$$y_i = \mathbf{w}^T \mathbf{h}(x_i) + \varepsilon_i \quad \varepsilon_i \sim N(0, \sigma^2)$$

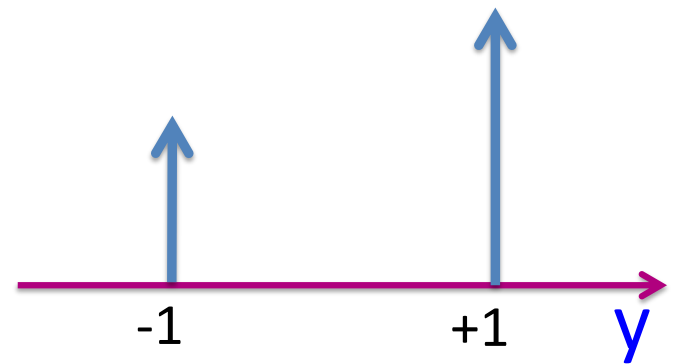
$$\rightarrow p(y | x, \mathbf{w}) = N(y; \mathbf{w}^T \mathbf{h}(x), \sigma^2)$$



Compare and contrast regression models

- Logistic regression

$$P(y|x, \mathbf{w}) = \begin{cases} \frac{1}{1 + e^{-\mathbf{w}^T h(x)}} & y = +1 \\ \frac{e^{-\mathbf{w}^T h(x)}}{1 + e^{-\mathbf{w}^T h(x)}} & y = -1 \end{cases}$$

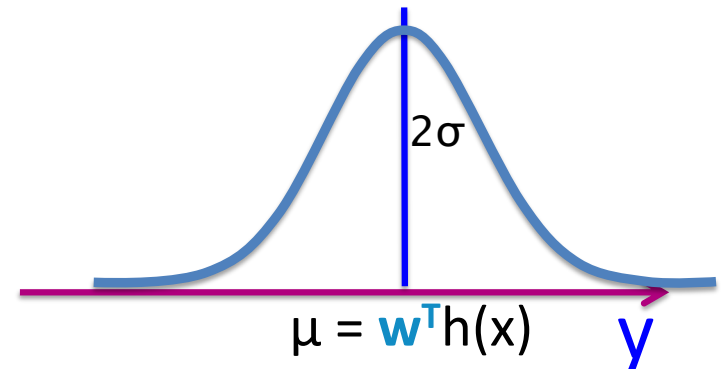


Compare and contrast regression models

- Linear regression with Gaussian errors

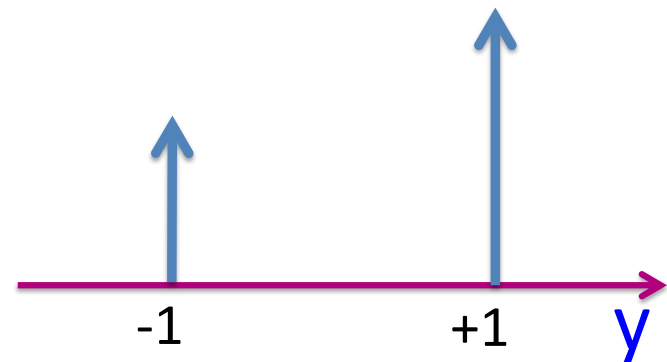
$$y_i = \mathbf{w}^T \mathbf{h}(x_i) + \varepsilon_i \quad \varepsilon_i \sim N(0, \sigma^2)$$

$$\rightarrow p(y | x, \mathbf{w}) = N(y; \mathbf{w}^T \mathbf{h}(x), \sigma^2)$$

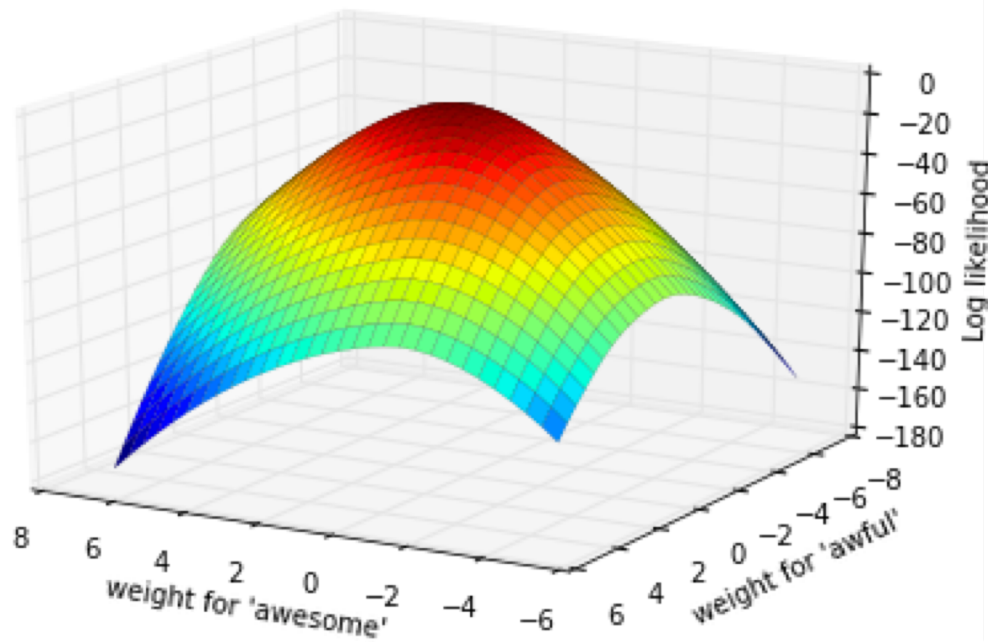


- Logistic regression

$$P(y | x, \mathbf{w}) = \begin{cases} \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{h}(x)}} & y = +1 \\ \frac{e^{-\mathbf{w}^T \mathbf{h}(x)}}{1 + e^{-\mathbf{w}^T \mathbf{h}(x)}} & y = -1 \end{cases}$$



Maximizing likelihood



No closed-form solution → use gradient ascent

Maximize function over all possible w_0, w_1, w_2

$$\max_{w_0, w_1, w_2} \prod_{i=1}^N P(y_i | \mathbf{x}_i, \mathbf{w})$$

$\ell(w_0, w_1, w_2)$ is a function of 3 variables

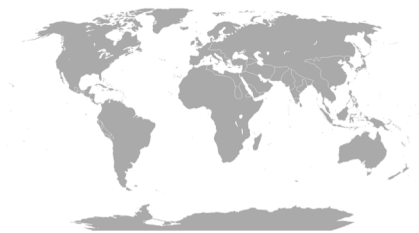
Encoding categorical inputs

Categorical inputs

- Numeric inputs:
 - #awesome, age, salary,...
 - Intuitive when multiplied by coefficient
 - e.g., 1.5 #awesome
- Categorical inputs:



Gender
(Male, Female,...)



Country of birth
(Argentina, Brazil, USA,...)



Zipcode
(10005, 98195,...)

Numeric value, but should be interpreted as category
(98195 not about 9x larger than 10005)

How do we multiply category by coefficient???
Must convert categorical inputs into numeric features

Encoding categories as numeric features

