

Models so far

We saw that a reasonable model for continuous outputs ($y \in \mathbb{R}$) is **linear regression**.

$$\text{Predict } y \in \text{ as } \begin{cases} y = \mathbf{x}^T \mathbf{w} & \text{(prediction function)} \\ p(y | \mathbf{x}, \mathbf{w}, \sigma^2) = \mathcal{N}(y | \mathbf{x}^T \mathbf{w}, \sigma^2) & \text{(probabilistic view)} \end{cases}$$

A reasonable model for *binary* outputs ($y \in \{0, 1\}$) is **logistic regression**:

$$\text{Predict } y \in \text{ as } \begin{cases} y = \mathbb{I}(\mathbf{x}^T \mathbf{w} > 0) & \text{(prediction function)} \\ p(y = 1 | \mathbf{x}, \mathbf{w}) = \sigma(\mathbf{x}^T \mathbf{w}) & \text{(probabilistic view)} \end{cases}$$

A reasonable model for *categorical* outputs ($y \in \{0, 1, \dots, C\}$) is **multinomial logistic regression**:

$$\text{Predict } y \in \text{ as } \begin{cases} y = \underset{c}{\operatorname{argmax}} \mathbf{x}^T \mathbf{w}_c & \text{(prediction function)} \\ p(y = c | \mathbf{x}, \mathbf{w}) = \operatorname{softmax}(\mathbf{x}^T \mathbf{W})_c & \text{(probabilistic view)} \end{cases}$$

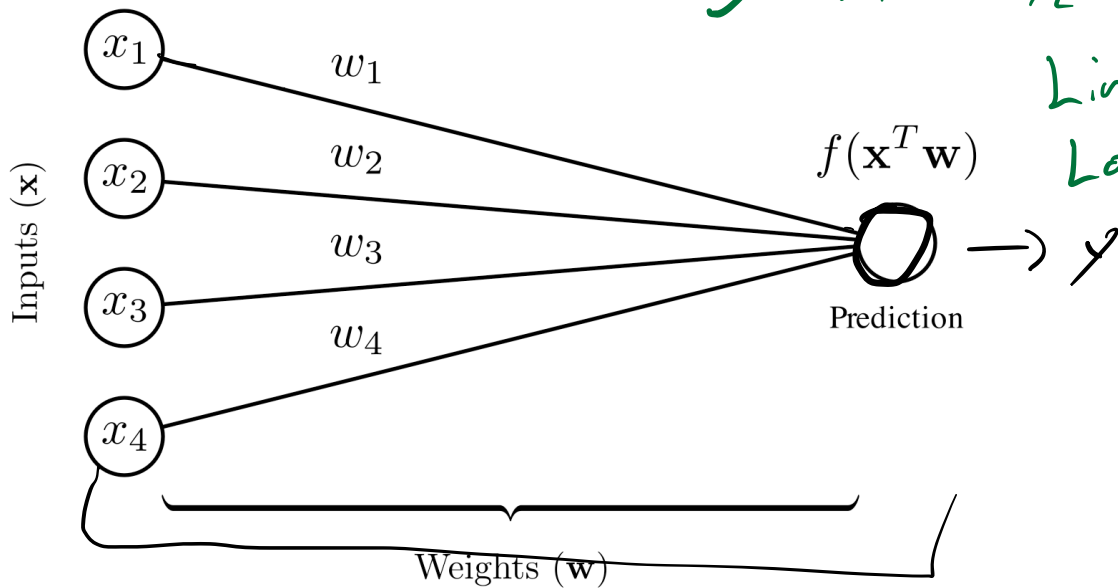
$$X = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ 1 \end{bmatrix} = \begin{bmatrix} \text{Weight} \\ \text{d13P} \\ \text{HP} \\ \vdots \end{bmatrix} \cdot \begin{bmatrix} w_1 \\ w_2 \\ \vdots \end{bmatrix} \quad \text{A new view}$$

Weights

$$f(x^T w) = x_1 w_1 + x_2 w_2 + x_3 w_3 + x_4 w_4$$

Linear $y = x^T w \quad f(a) = a$

Logistic $y = \mathbb{I}(x^T w > 0) \quad f(a) = \mathbb{I}(a > 0)$



Linear Predictions

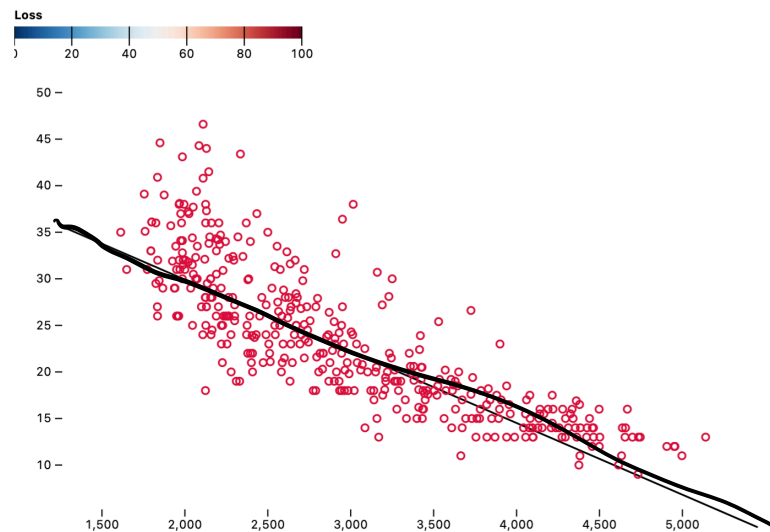
$$f(\mathbf{x}) = \mathbf{x}^T \mathbf{w} = \sum_{i=1}^n x_i w_i$$

Handwritten: $\begin{bmatrix} x \\ 1 \end{bmatrix} \begin{bmatrix} w \\ b \end{bmatrix}$

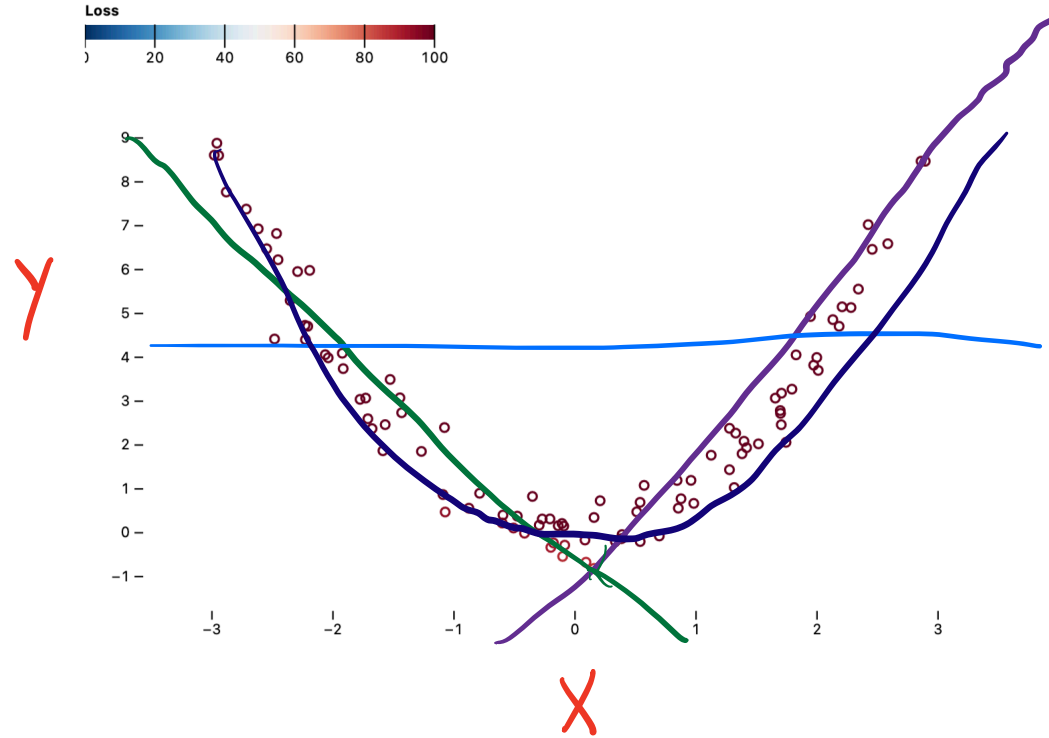
Predicted MPG = $f(\mathbf{x}) =$

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

Handwritten: \hookrightarrow



Non-Linear Data

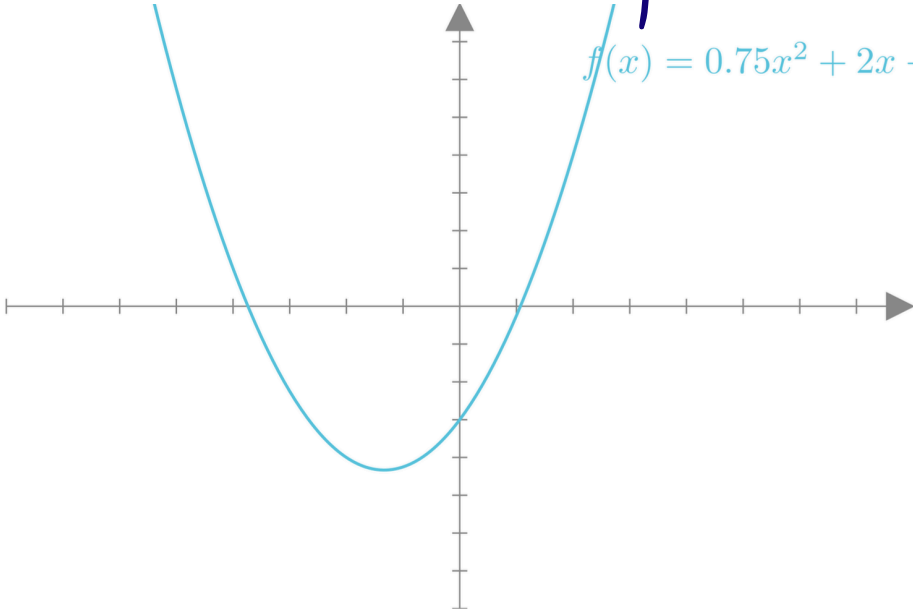


Polynomial functions

Quadratic functions

$$f(x) = w_2x^2 + w_1x + b$$

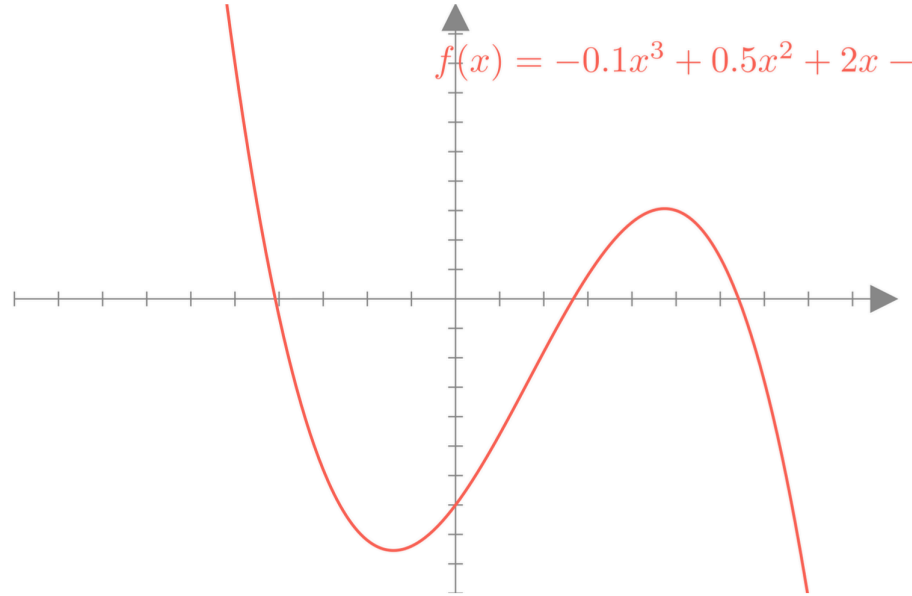
$$f(x) = 0.75x^2 + 2x - 3$$



Cubic functions

$$f(x) = w_3x^3 + w_2x^2 + w_1x + b$$

$$f(x) = -0.1x^3 + 0.5x^2 + 2x - 7$$



Polynomial functions of multiple inputs

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

artificial

$$y = \underbrace{f(\mathbf{x}) = \underbrace{w_5 x_2^2 + w_4 x_1^2 + w_3 x_1 x_2 + w_2 x_2 + w_1 x_1 + b}}_{\text{artificial}}$$

degree 4 :

$$\left. \begin{aligned} f(x) &= w_1 x_1^4 + b \\ f(x) &= w_1 x_1^2 x_2^2 \\ &\quad \vee x_1^3 x_2 \end{aligned} \right\} \text{degree 4}$$

The degree of a polynomial is the largest exponent in any term of the polynomial

Polynomial functions as vector functions

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

$$f(\mathbf{x}) = w_5 x_2^2 + w_4 x_1^2 + w_3 x_1 x_2 + w_2 x_2 + w_1 x_1 + b$$

$$w_5 x_2^2 + w_4 x_1^2 + w_3 x_1 x_2 + w_2 x_2 + w_1 x_1 + b =$$

$$\begin{bmatrix} x_1 \\ x_2 \\ x_1 x_2 \\ x_1^2 \\ x_2^2 \\ 1 \end{bmatrix} \cdot \begin{bmatrix} w_1 \\ w_2 \\ w_3 \\ w_4 \\ w_5 \\ b \end{bmatrix}$$

$\hat{\phi}(\mathbf{x})$: Feature transform of \mathbf{x}

$$f(\mathbf{x}) = \phi(\mathbf{x})^T \mathbf{w}$$

Quadratic feature transforms

Let's consider the mapping from \mathbf{x} to powers of the elements of \mathbf{x} . We'll call this mapping ϕ :

$$\underbrace{\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}}_{\mathbf{x}} \xrightarrow{\phi} \begin{bmatrix} x_1 \\ x_2 \\ x_1 x_2 \\ x_1^2 \\ x_2^2 \\ 1 \end{bmatrix}$$

$$\phi(\mathbf{x}) = \begin{bmatrix} x_1 \\ x_2 \\ x_1 x_2 \\ x_1^2 \\ x_2^2 \\ 1 \end{bmatrix}$$

$$f(\mathbf{x}) = \phi(\mathbf{x})^T \mathbf{w}$$

$$f(\mathbf{x}) = \phi(\mathbf{x})^T \mathbf{w} = w_2 x_1^2 + w_1 x_1 + b, \quad \phi(\mathbf{x}) = \begin{bmatrix} x_1 \\ x_1^2 \\ 1 \end{bmatrix}$$

Fitting quadratic regression

Prediction function

$$\underbrace{f(\mathbf{x}) = \phi(\mathbf{x})^T \mathbf{w}}_{\text{predict } y}, \quad \phi(\mathbf{x}) = \begin{bmatrix} x_1 \\ x_2 \\ x_1 x_2 \\ x_1^2 \\ x_2^2 \\ 1 \end{bmatrix}$$

Probabilistic model

$$y_i \sim \mathcal{N}(\underbrace{\phi(\mathbf{x}_i)^T \mathbf{w}}_{\text{prediction}}, \sigma^2)$$

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

Negative log-likelihood loss

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

$$= \frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \phi(\mathbf{x}_i)^T \mathbf{w})^2 + N \log \sigma \sqrt{2\pi}$$

$$\mathbf{w}^{(i+1)} \leftarrow \mathbf{w}^{(i)} - \alpha \nabla_{\mathbf{w}} \text{loss}(\mathbf{x}, y, \mathbf{w})$$

prediction func

$$\mathbf{x}^T \mathbf{w} \rightarrow \phi(\mathbf{x})^T \mathbf{w}$$

Fitting quadratic regression

Prediction function

$$f(\mathbf{x}) = \phi(\mathbf{x})^T \mathbf{w}, \quad \phi(\mathbf{x}) = \begin{bmatrix} x_1 \\ x_2 \\ x_1 x_2 \\ x_1^2 \\ x_2^2 \\ 1 \end{bmatrix}$$

Probabilistic model

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

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

Negative log-likelihood loss

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

$$= \frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \phi(\mathbf{x}_i)^T \mathbf{w})^2$$

We see that the gradient doesn't change, it simply involves $\phi(\mathbf{x}_i)$ instead of \mathbf{x}_i .

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

Quadratic logistic regression

Prediction function

$$f(\mathbf{x}) = \mathbb{I}(\phi(\mathbf{x})^T \mathbf{w} \geq 0), \quad \phi(\mathbf{x}) = \begin{bmatrix} x_1 \\ x_2 \\ x_1 x_2 \\ x_1^2 \\ x_2^2 \\ 1 \end{bmatrix}$$

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

Probabilistic model

$p(y_i = 1 | \mathbf{x}_i, \mathbf{w})$

$$y_i \sim \text{Bernoulli}(\sigma(\phi(\mathbf{x}_i)^T \mathbf{w})), \quad p(y_i = 1 | \mathbf{x}_i, \mathbf{w}) = \sigma(\phi(\mathbf{x}_i)^T \mathbf{w})$$

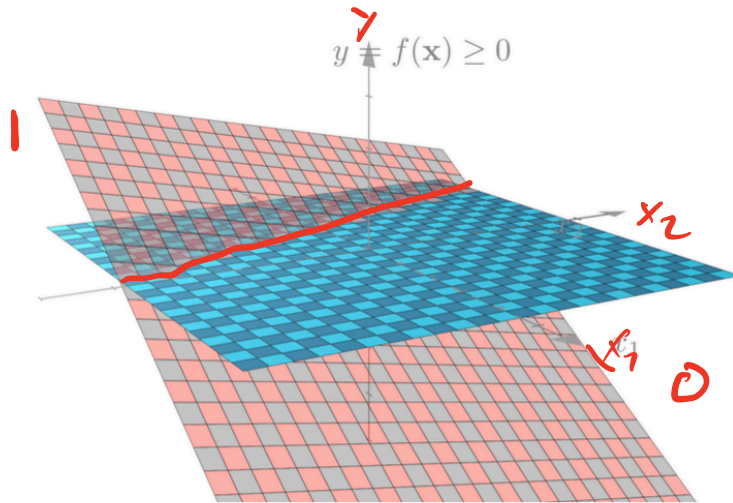
$$f(x) = \mathbb{I}(x_1 w_1 + x_2 w_2 + x_1 x_2 w_3 + x_1^2 w_4 + x_2^2 w_5 + b \geq 0)$$

Negative log-likelihood loss

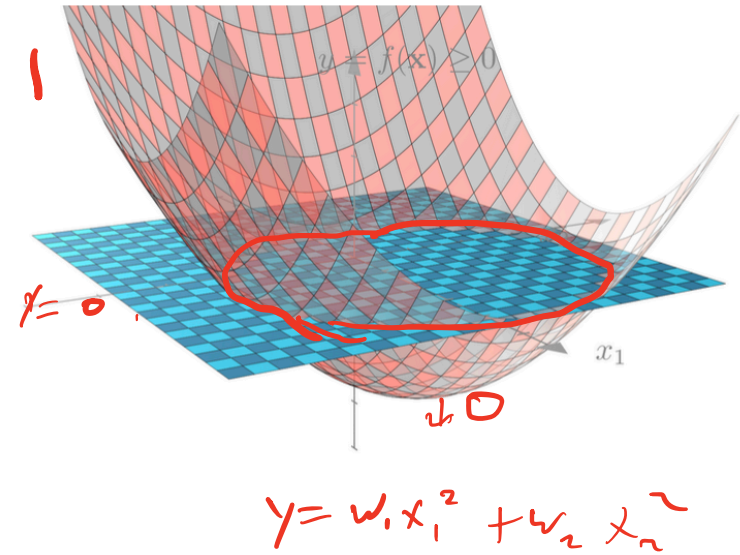
$$\text{NLL}(\mathbf{w}, \mathbf{X}, \mathbf{y}) = - \sum_{i=1}^N \log \sigma((2y_i - 1) \phi(\mathbf{x}_i)^T \mathbf{w})$$

Quadratic decision boundary

Linear decision boundary



Quadratic decision boundary



Cubic feature transform

$$\underset{\uparrow}{[x_1]} \xrightarrow{\phi} \underset{\uparrow}{\begin{bmatrix} x_1 \\ x_1^2 \\ x_1^3 \\ 1 \end{bmatrix}}$$

Prediction function (regression)

$$f(\mathbf{x}) = \phi(\mathbf{x})^T \mathbf{w} = w_3 x_1^3 + w_2 x_1^2 + w_1 x_1 + b, \quad \phi(\mathbf{x}) = \begin{bmatrix} x_1 \\ x_1^2 \\ x_1^3 \\ 1 \end{bmatrix}$$

General feature transforms

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

$$\Rightarrow \begin{bmatrix} x_1 \\ x_2 \\ x_1^2 \\ x_2^2 \\ x_1 x_2 \\ \sin(x_1) \\ \sin(x_2) \end{bmatrix}$$

$$f(\mathbf{x}) = \phi(\mathbf{x})^T \mathbf{w}, \quad \phi(\mathbf{x}) =$$

$$f(\mathbf{x}) = \underbrace{w_1 x_1 + w_2 x_2 + w_3 x_1^2 + w_4 x_2^2 + w_5 x_1 x_2 + w_6 \sin(x_1) + w_7 \sin(x_2)}$$

General feature transforms

$$X = N \times 1 \quad \begin{bmatrix} x_1 & x_2 \\ x_1^2 & x_2^2 \end{bmatrix} X$$

$$\begin{bmatrix} x_1^2 & x_2^2 \\ x_1 x_2 & \sin(x_1) \end{bmatrix}$$

$$f(\mathbf{x}) = \phi(\mathbf{x})^T \mathbf{w}, \quad \phi(\mathbf{x}) =$$

$$\begin{bmatrix} x_1 \\ x_2 \\ x_1^2 \\ x_2^2 \\ x_1 x_2 \\ \sin(x_1) \\ \sin(x_2) \end{bmatrix}$$

```
squared_X = X ** 2 # x^2
cross_X = X[:, :1] * X[:, 1:2] # x_1 * x_2
sin_X = np.sin(X) # sin(x)
transformedX = np.concatenate([X, squared_X, cross_X, sin_X], axis=1)
```

$$\Phi X = \begin{bmatrix} x_1 & x_2 \\ x_1^2 & x_2^2 \end{bmatrix} \begin{bmatrix} x_1 x_2 \\ \sin(x_1) \end{bmatrix}$$

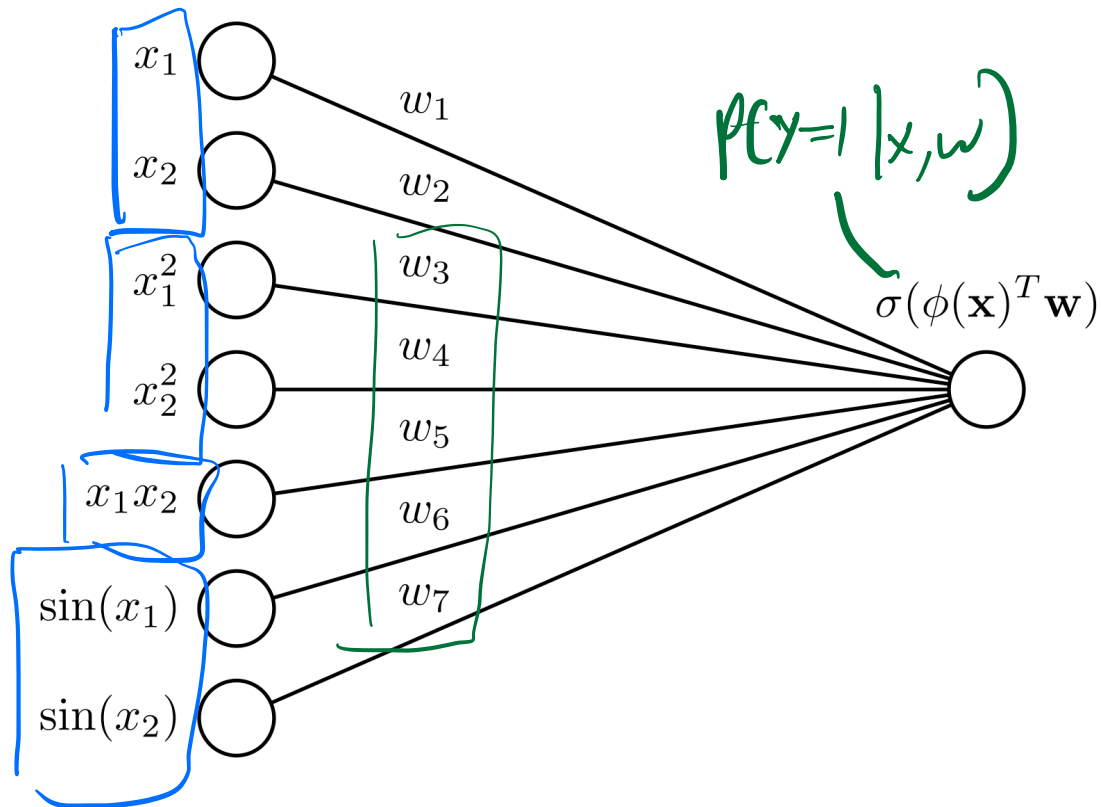
$$f(\mathbf{x}) = w_1 x_1 + w_2 x_2 + w_3 x_1^2 + w_4 x_2^2 + w_5 x_1 x_2 + w_6 \sin(x_1) + w_7 \sin(x_2)$$

$$f(\mathbf{x} | \mathbf{w})$$

Non-linear logistic regression

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

$$f(\mathbf{x}) = \phi(\mathbf{x})^T \mathbf{w}, \quad \phi(\mathbf{x}) = \begin{bmatrix} x_1 \\ x_2 \\ x_1^2 \\ x_2^2 \\ x_1 x_2 \\ \sin(x_1) \\ \sin(x_2) \end{bmatrix}$$



Learning linear functions

We've already seen that we can *learn* a function by defining our function in terms of a set of *parameters* \mathbf{w} :

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

and then minimizing a *loss* as a function of \mathbf{w}

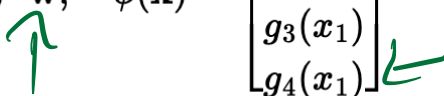
$$\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmin}} \operatorname{Loss}(\mathbf{w})$$

Which we can do with gradient descent:

$$\mathbf{w}^{(k+1)} \leftarrow \mathbf{w}^{(k)} - \alpha \nabla_{\mathbf{w}} \operatorname{Loss}(\mathbf{w})$$

Can we learn the functions in our feature transform?

$$\begin{aligned} g_1(x) &= x^2 \\ g_2(x) &= x \\ g_3(x) &= \sin(x) \end{aligned}$$

$$f(\mathbf{x}) = \phi(\mathbf{x})^T \mathbf{w}, \quad \phi(\mathbf{x}) = \begin{bmatrix} g_1(x_1) \\ g_2(x_1) \\ g_3(x_1) \\ g_4(x_1) \end{bmatrix}$$


Can we learn the functions in our feature transform?

$$g_i(\mathbf{x}) = \sigma(\mathbf{x}^T \mathbf{w}_i)$$

$$f(\mathbf{x}) = \phi(\mathbf{x})^T \mathbf{w}, \quad \phi(\mathbf{x}) = \begin{bmatrix} g_1(x_1) \\ g_2(x_1) \\ g_3(x_1) \\ g_4(x_1) \end{bmatrix}$$

Logistic Regression

$$f(\mathbf{x}) = \phi(\mathbf{x})^T \mathbf{w}_0, \quad \phi(\mathbf{x}) = \begin{bmatrix} \sigma(\mathbf{x}^T \mathbf{w}_1) \\ \sigma(\mathbf{x}^T \mathbf{w}_2) \\ \sigma(\mathbf{x}^T \mathbf{w}_3) \\ \sigma(\mathbf{x}^T \mathbf{w}_4) \end{bmatrix}$$

parameters of g_1

parameters of g_2

parameters of prediction func.

All parameters : w_0, w_1, w_2, w_3, w_4

A simple learned transform

Let's look at a very simple example:

$$w_i = \begin{bmatrix} w_{i1} \\ w_{i2} \end{bmatrix}$$

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, \quad \mathbf{w}_0 = \begin{bmatrix} w_{01} \\ w_{02} \end{bmatrix}$$
$$\underbrace{f(\mathbf{x}) = \phi(\mathbf{x})^T \mathbf{w}_0}_{\text{prediction function}} \quad \phi(\mathbf{x}) = \begin{bmatrix} \sigma(\mathbf{x}^T \mathbf{w}_1) \\ \sigma(\mathbf{x}^T \mathbf{w}_2) \\ \sigma(\mathbf{x}^T \mathbf{w}_3) \end{bmatrix} = \begin{bmatrix} \sigma(x_1 w_{11} + x_2 w_{12}) \\ \sigma(x_1 w_{21} + x_2 w_{22}) \\ \sigma(x_1 w_{31} + x_2 w_{32}) \end{bmatrix}$$

Handwritten notes: A bracket labeled w_i groups the w_{i1} and w_{i2} terms in the $\phi(\mathbf{x})$ vector. A bracket labeled 2×1 groups the two input terms $x_1 w_{i1}$ and $x_2 w_{i2}$ in each $\sigma(\cdot)$ function.

In this case, we can write out our prediction function explicitly as:

$$f(x) = \phi(x)_1 w_{01} + \phi(x)_2 w_{02} + \phi(x)_3 w_{03}$$

$$= \underbrace{\sigma(x_1 w_{11} + x_2 w_{12})}_{\text{activation}} w_{01} + \sigma(x_1 w_{21} + x_2 w_{22}) w_{02} + \sigma(x_1 w_{31} + x_2 w_{32}) w_{03}$$

A simple learned transform

Let's look at a very simple example:

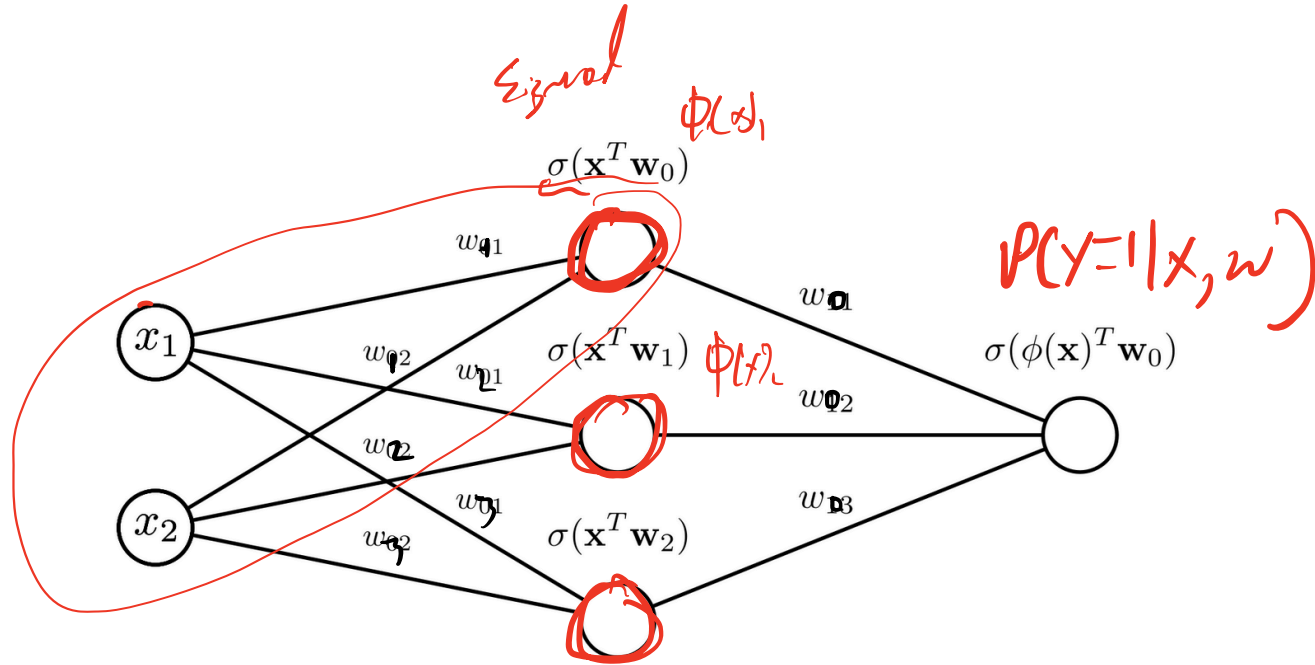
$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, \quad \mathbf{w}_0 = \begin{bmatrix} w_{01} \\ w_{02} \end{bmatrix}$$

$$f(\mathbf{x}) = \phi(\mathbf{x})^T \mathbf{w}_0, \quad \phi(\mathbf{x}) = \begin{bmatrix} \sigma(\mathbf{x}^T \mathbf{w}_1) \\ \sigma(\mathbf{x}^T \mathbf{w}_2) \\ \sigma(\mathbf{x}^T \mathbf{w}_3) \end{bmatrix} = \begin{bmatrix} \sigma(x_1 w_{11} + x_2 w_{12}) \\ \sigma(x_1 w_{21} + x_2 w_{22}) \\ \sigma(x_1 w_{31} + x_2 w_{32}) \end{bmatrix}$$

In this case, we can write out our prediction function explicitly as:

$$f(\mathbf{x}) = w_{01} \cdot \sigma(x_1 w_{11} + x_2 w_{12}) + w_{02} \cdot \sigma(x_1 w_{21} + x_2 w_{22}) + w_{03} \cdot \sigma(x_1 w_{31} + x_2 w_{32})$$

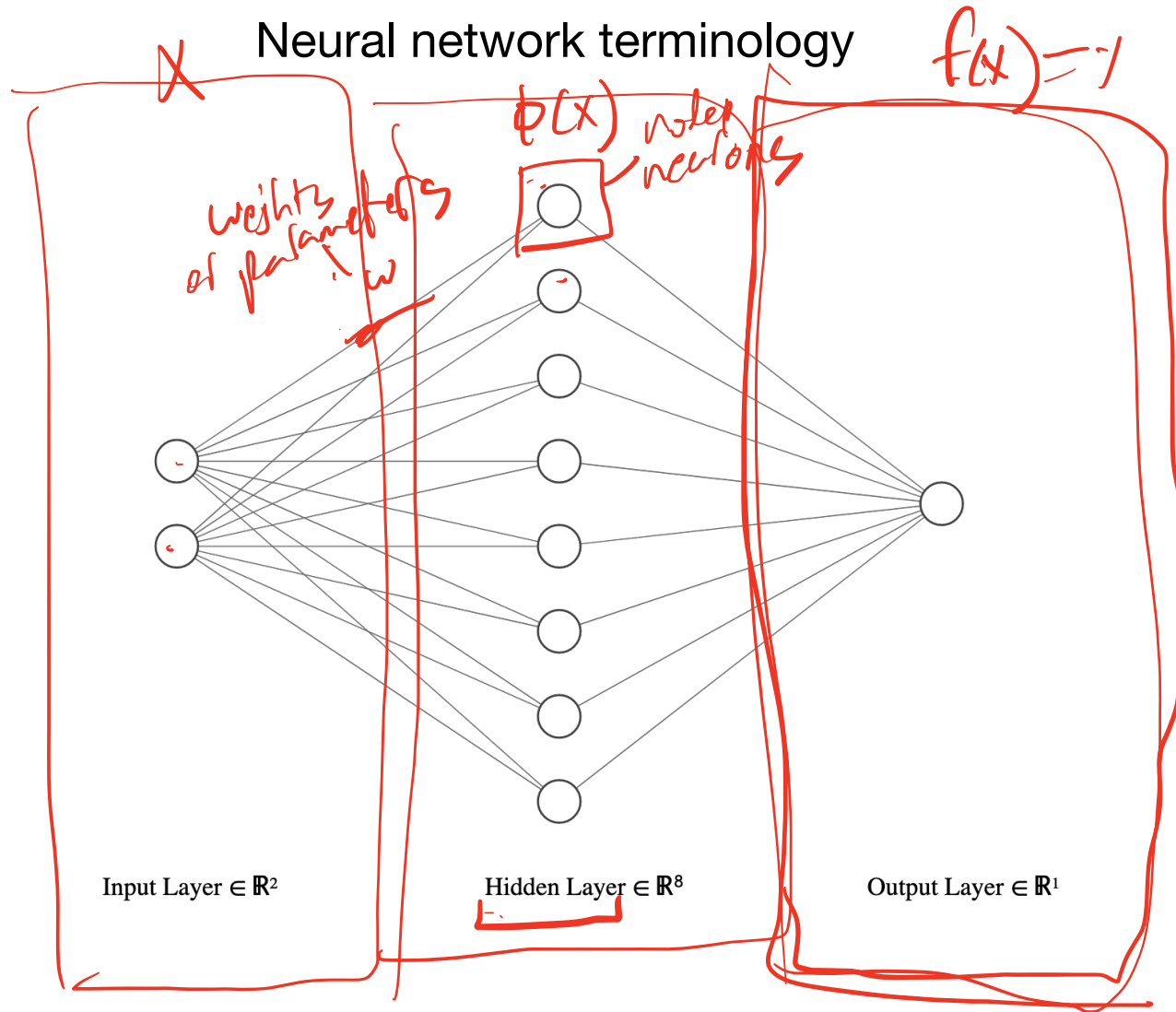
A neural network!



In this case, we can write out our prediction function explicitly as:

$$f(\mathbf{x}) = w_{01} \cdot \sigma(x_1 w_{11} + x_2 w_{12}) + w_{02} \cdot \sigma(x_1 w_{21} + x_2 w_{22}) + w_{03} \cdot \sigma(x_1 w_{31} + x_2 w_{32})$$

Neural network terminology



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

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

-