

Feature transforms

Story so far

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

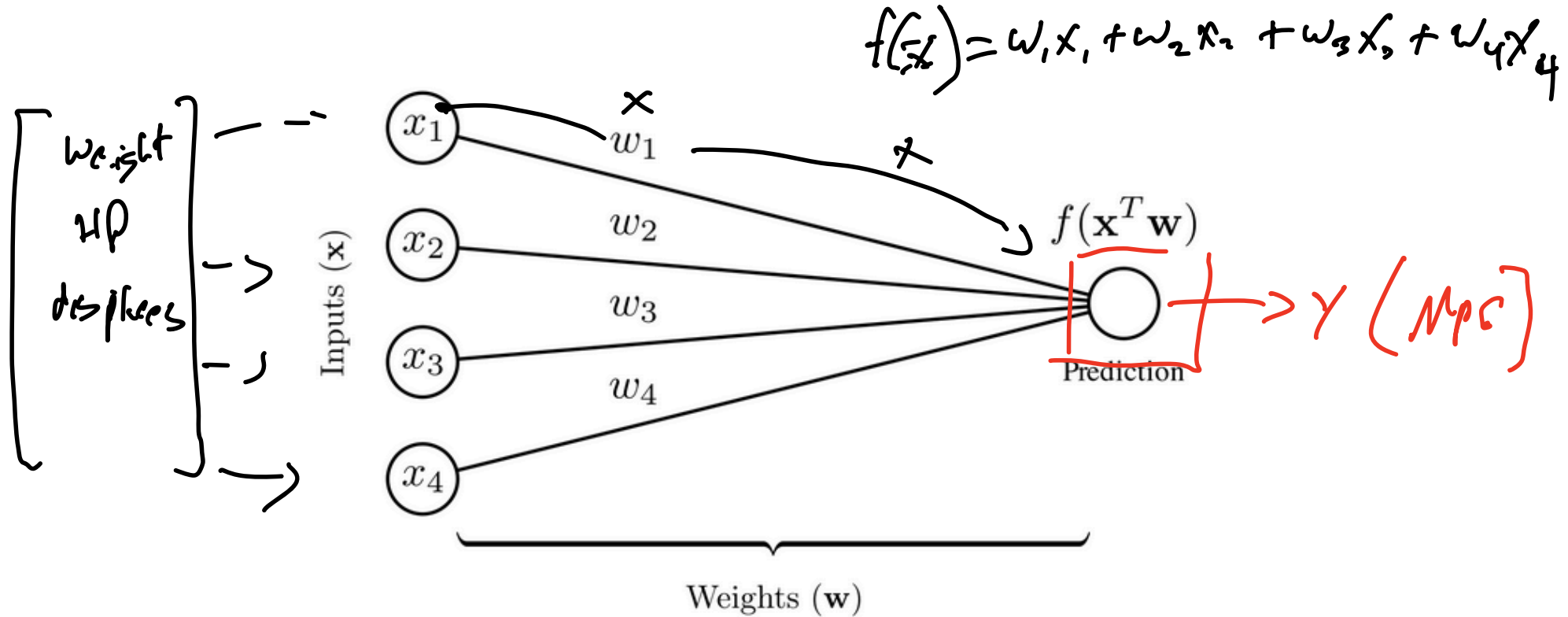
A reasonable model for *binary* outputs ($y \in \{0, 1\}$) is **logistic regression**:

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

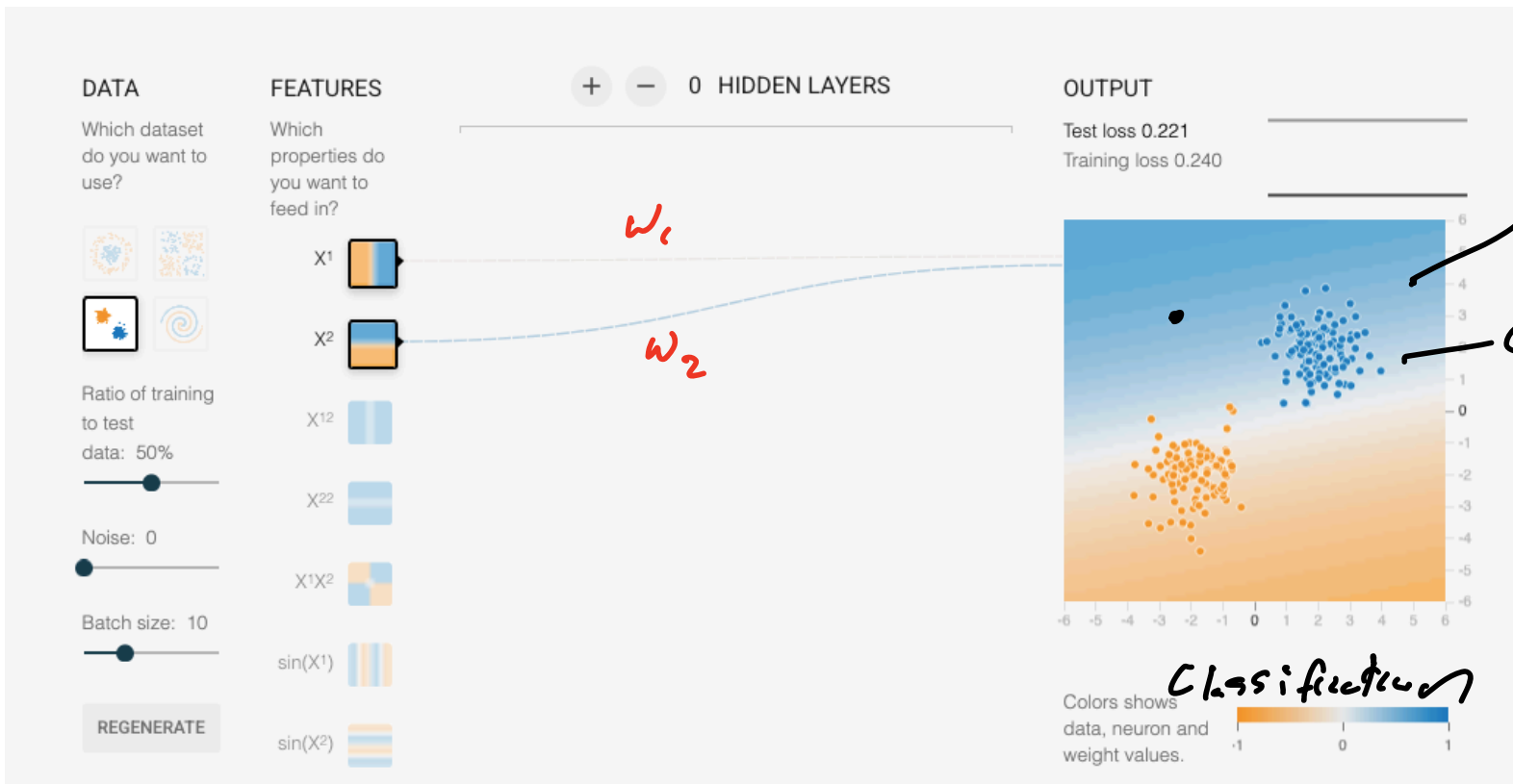
A reasonable model for *categorical* outputs ($y \in \{0, 1, \dots, C\}$) is **multinomial logistic regression**:

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

A new visualization

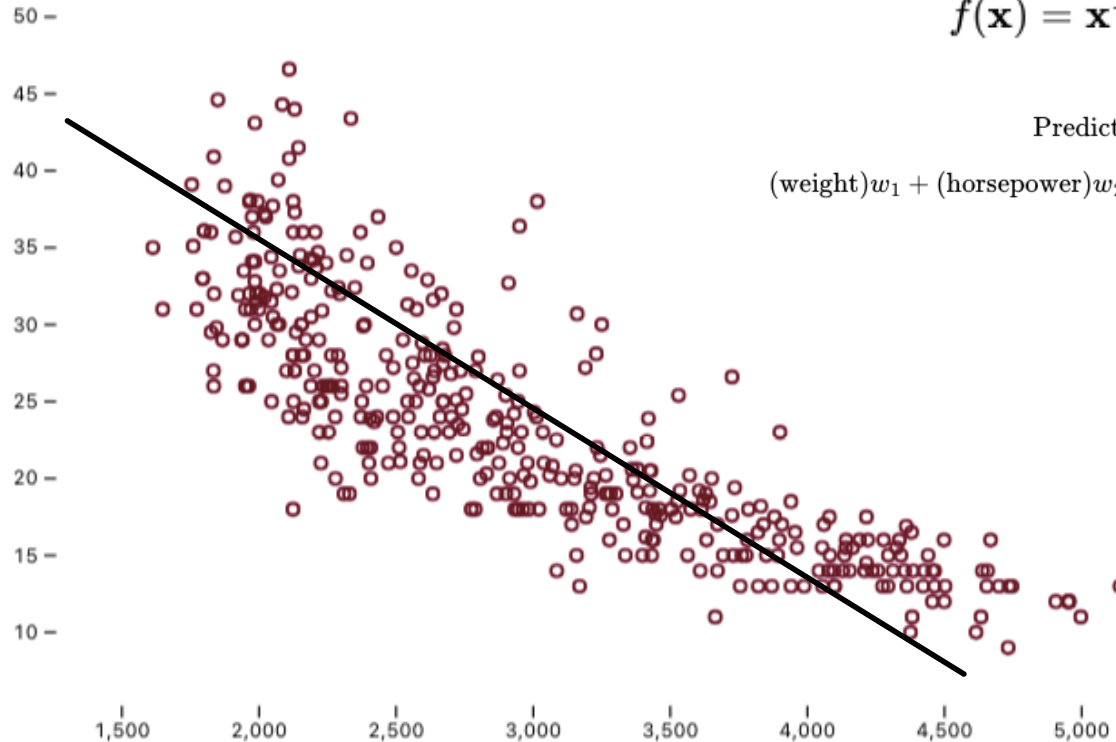


playground.tensorflow.org



(Approx.) Linear data

mpg ↑



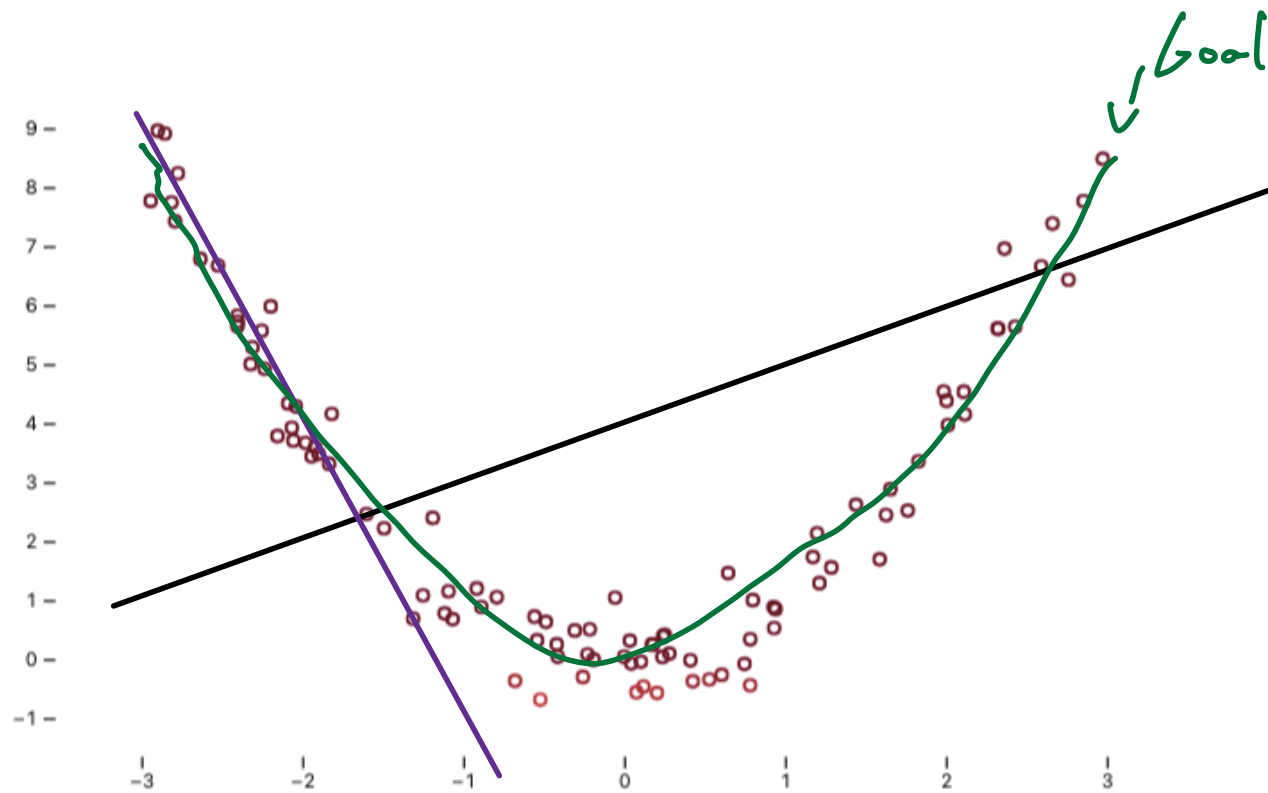
$$f(\mathbf{x}) = \mathbf{x}^T \mathbf{w} = \sum_{i=1}^n x_i w_i$$

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

(weight) w_1 + (horsepower) w_2 + (displacement) w_3 + (0-60mph) w_4 + b

weights →

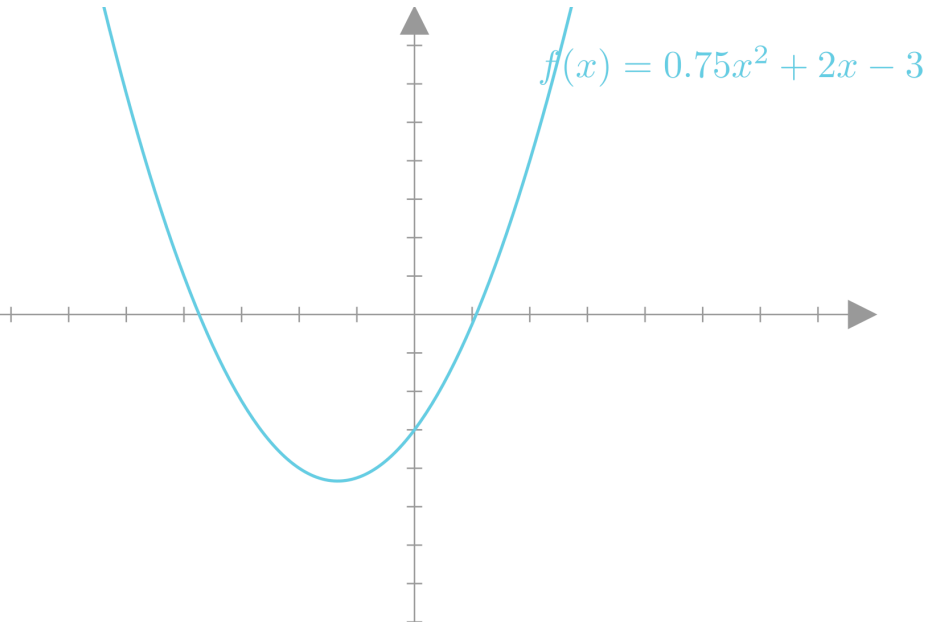
Non-Linear data



Polynomial functions

Quadratic function

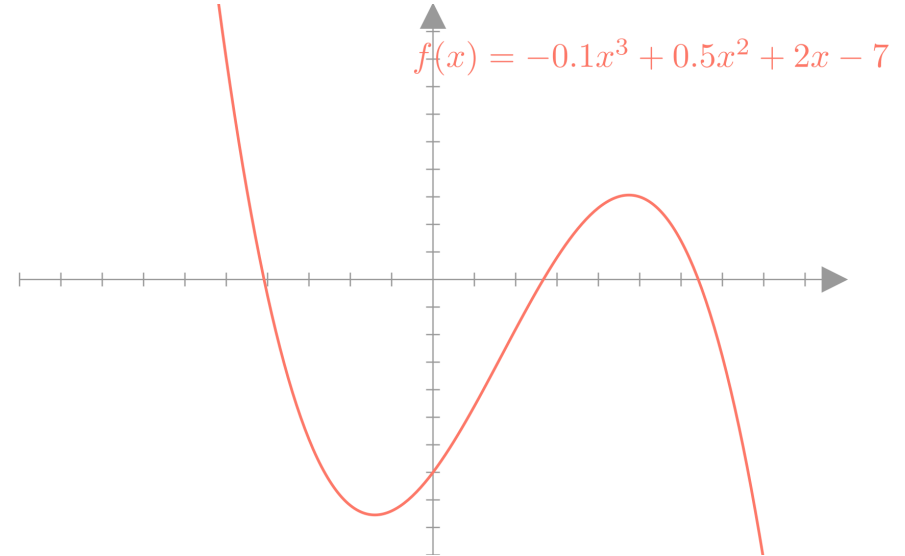
$$f(x) = w_2x^2 + w_1x + b$$



Degree (highest power): 2

Cubic function

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



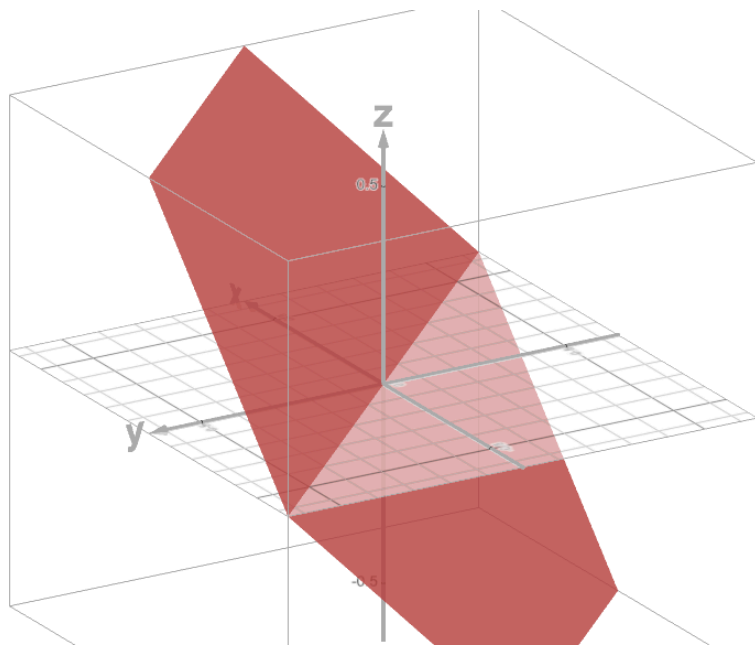
Degree (highest power): 3

Polynomial functions of multiple inputs

$$f(x, y) = w_5x^2 + w_4y^2 + w_3xy + w_2x + w_1y + b$$

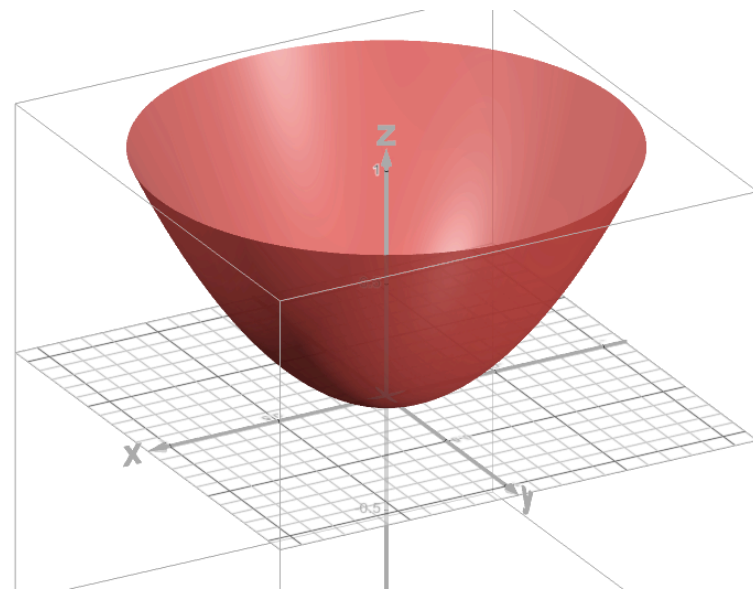
$$f(\mathbf{x}) = w_5x_2^2 + w_4x_1^2 + w_3x_1x_2 + w_2x_2 + w_1x_1 + b$$

Linear function of 2-inputs (plane)



Quadratic function of 2-inputs (paraboloid)

$$f(\mathbf{x}) = w_5x_2^2 + w_4x_1^2 + w_3x_1x_2 + w_2x_2 + w_1x_1 + b$$



Degree of a polynomial function

Largest (total) exponent in any term

Degree 2 polynomials

$$f(x, y) = w_5x^2 + w_4y^2 + w_3xy + w_2x + w_1y + b$$

Degree 4 polynomials

$$f(x, y) = 3x^4 + 2xy + y - 2$$

$$f(x, y) = -2x^2y^2 + 2x^3 + y^2 - 5$$

Polynomial functions as vector operations

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

$$f(x) = w + b = \begin{bmatrix} x \\ 1 \end{bmatrix} \cdot \begin{bmatrix} w \\ b \end{bmatrix}$$

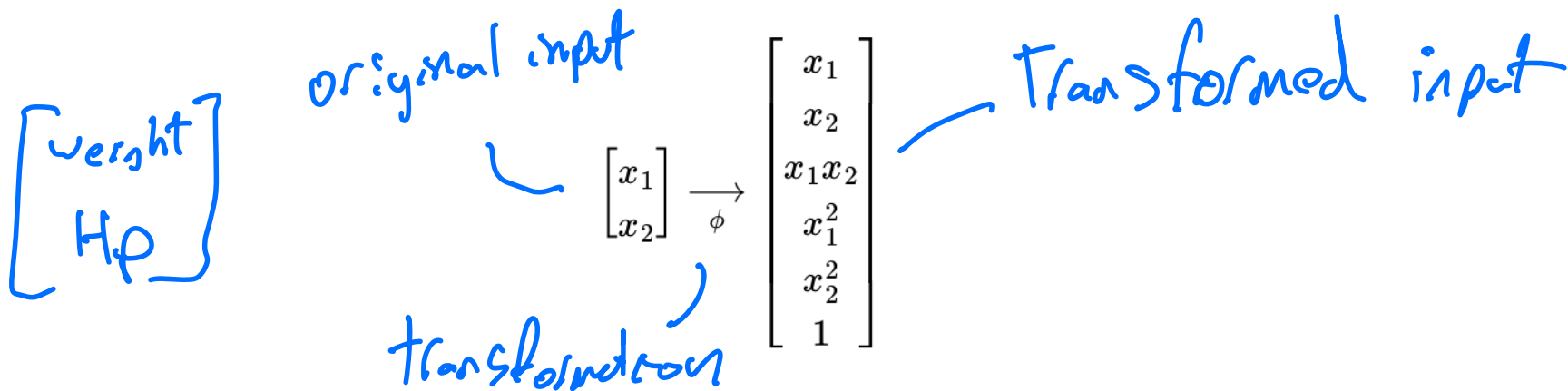
powers of
inputs

parameters

$$f(\tilde{x}) = w_5 x_2^2 + w_4 x_1^2 + w_3 x_1 x_2 + w_2 x_2 + w_1 x_1 + b = \begin{bmatrix} x_1 \\ x_2 \\ x_1 x_2 \\ x_1^2 \\ x_2^2 \\ 1 \end{bmatrix} \cdot \begin{bmatrix} w_1 \\ w_2 \\ w_3 \\ w_4 \\ w_5 \\ b \end{bmatrix}$$

Polynomial functions as vector operations

$$w_5x_2^2 + w_4x_1^2 + w_3x_1x_2 + w_2x + w_1y + b = \begin{bmatrix} x_1 \\ x_2 \\ x_1x_2 \\ x_1^2 \\ x_2^2 \\ 1 \end{bmatrix} \cdot \begin{bmatrix} w_1 \\ w_2 \\ w_3 \\ w_4 \\ w_5 \\ b \end{bmatrix}$$



Polynomial functions as vector operations

$$f: \mathbb{R}^2 \rightarrow \mathbb{R}^6 \quad f(\mathbf{x}) = w_1 \underline{x_1^2} + w_2 \underline{x_2^2} + w_3 \underline{x_1 x_2} + v_4 \underline{x_1} + v_5 \underline{x_2} + b$$

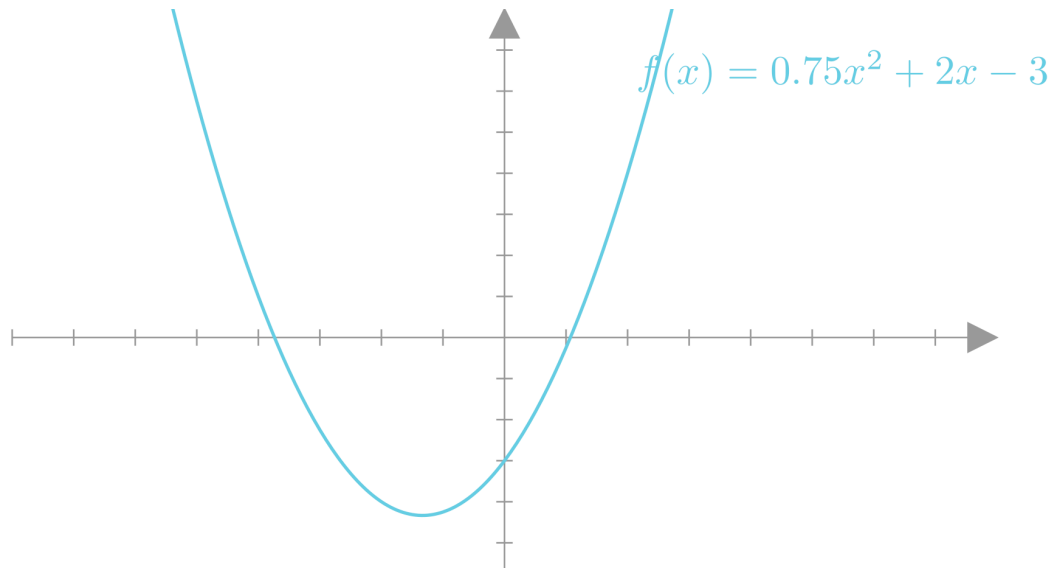
$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \xrightarrow{\phi} \begin{bmatrix} x_1 \\ x_2 \\ x_1 x_2 \\ x_1^2 \\ x_2^2 \\ 1 \end{bmatrix}$$

$$\phi(\mathbf{x}) = \begin{bmatrix} x_1 \\ x_2 \\ x_1 x_2 \\ x_1^2 \\ x_2^2 \\ 1 \end{bmatrix} \quad \begin{bmatrix} w_3 \\ w_4 \\ w_3 \\ w_2 \\ v_1 \\ b \end{bmatrix}$$

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

Quadratic transformation func.

Quadratic function as a feature transform



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

Fitting quadratic regression

$$x^T w \rightarrow \phi(x)^T w$$

Prediction function

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

Negative log-likelihood loss

$$\begin{aligned} \text{Loss}(\mathbf{w}) &= \text{NLL}(\mathbf{w}, \mathbf{X}, \mathbf{y}) = - \sum_{i=1}^N \log p(y_i \mid \mathbf{x}_i, \mathbf{w}) \\ &= \frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \phi(\mathbf{x}_i)^T \mathbf{w})^2 + N \log \sigma \sqrt{2\pi} \end{aligned}$$

Probabilistic model

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

Optimization problem

$$\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmin}} \text{NLL}(\mathbf{w}, \mathbf{X}, \mathbf{y})$$

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

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

Negative log-likelihood loss

$$\begin{aligned} \text{Loss}(\mathbf{w}) &= \text{NLL}(\mathbf{w}, \mathbf{X}, \mathbf{y}) = - \sum_{i=1}^N \log p(y_i \mid \mathbf{x}_i, \mathbf{w}) \\ &= \frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \phi(\mathbf{x}_i)^T \mathbf{w})^2 + N \log \sigma \sqrt{2\pi} \end{aligned}$$

What is the gradient of the log-likelihood with respect to \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

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

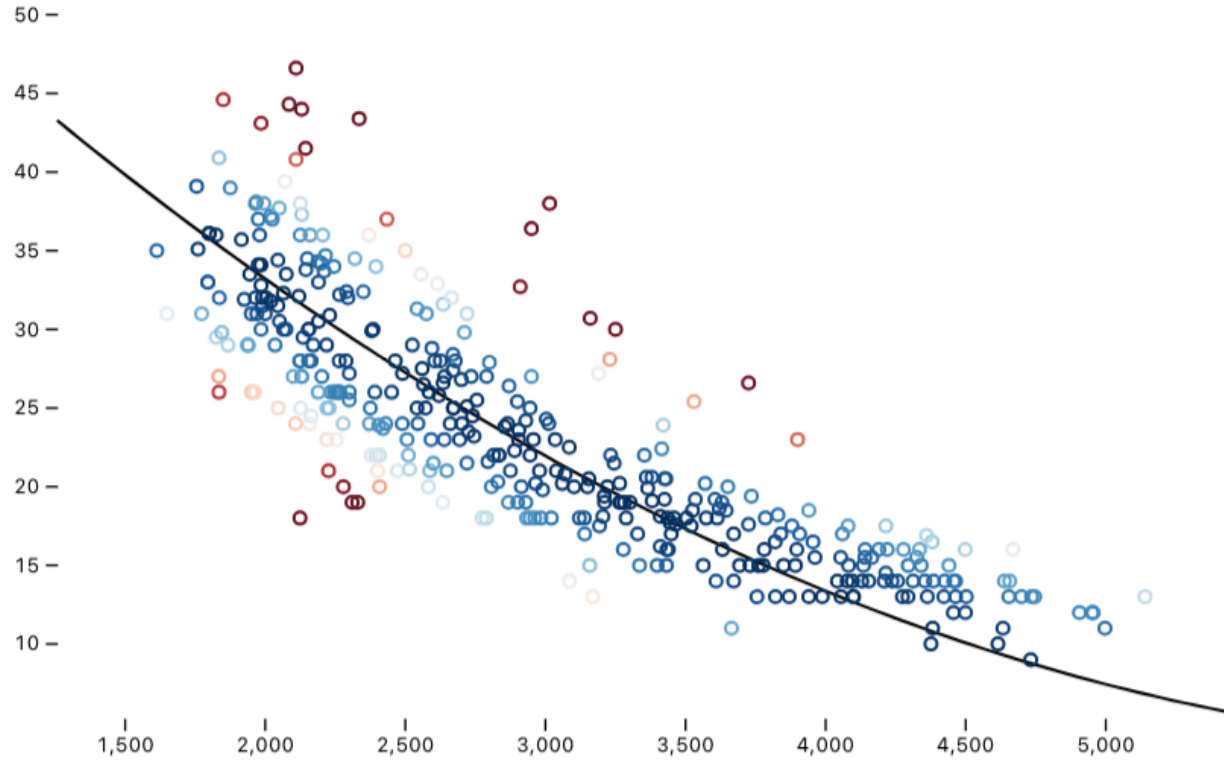
Negative log-likelihood loss

$$\begin{aligned} \text{Loss}(\mathbf{w}) &= \text{NLL}(\mathbf{w}, \mathbf{X}, \mathbf{y}) = - \sum_{i=1}^N \log p(y_i \mid \mathbf{x}_i, \mathbf{w}) \\ &= \frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \phi(\mathbf{x}_i)^T \mathbf{w})^2 + N \log \sigma \sqrt{2\pi} \end{aligned}$$

What is the gradient of the log-likelihood with respect to \mathbf{w} ?

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

Quadratic regression on real data



Quadratic logistic regression

Prediction function

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

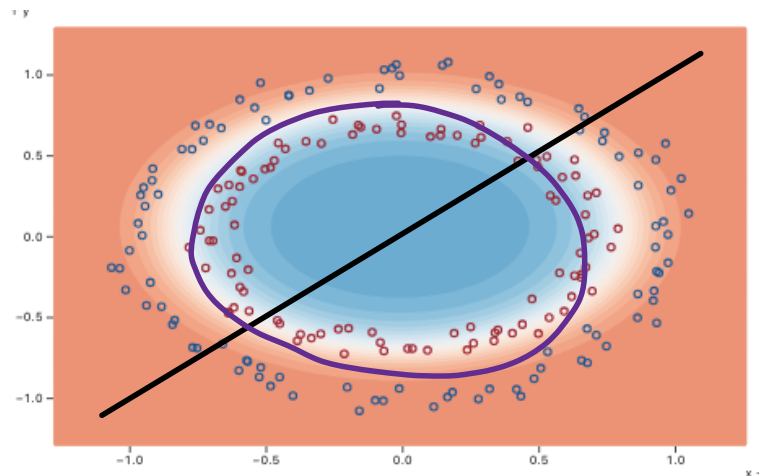
Probabilistic model

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

Negative log-likelihood loss

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

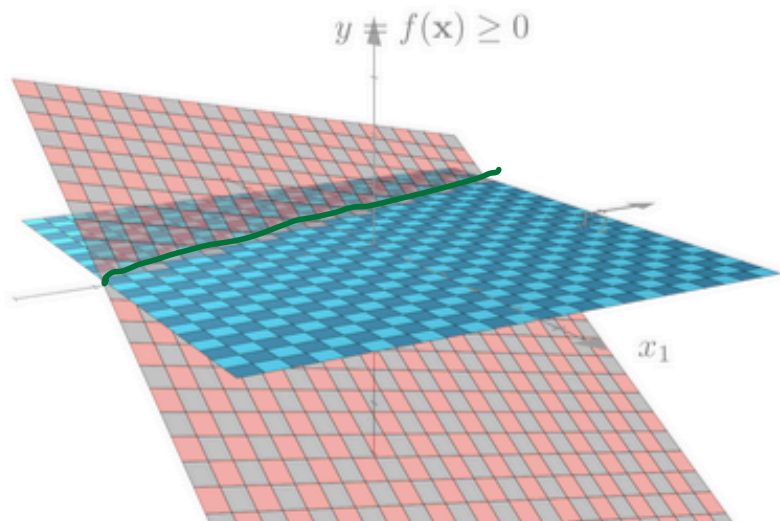
With two inputs



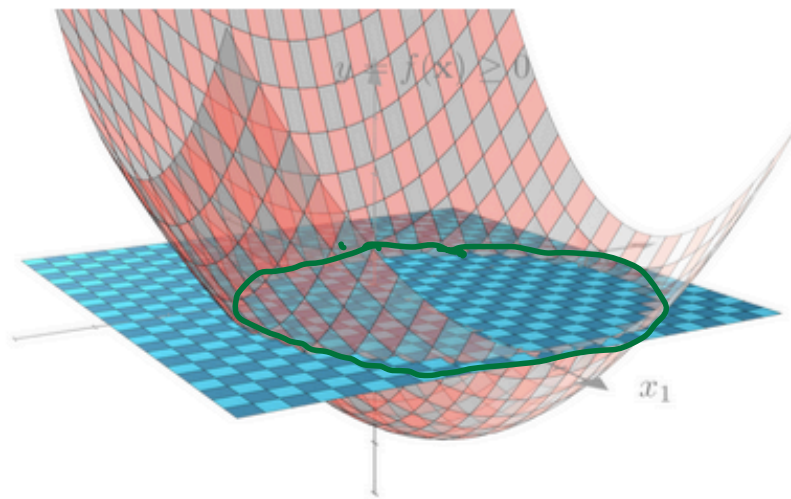
Quadratic decision boundaries

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

Linear decision boundary

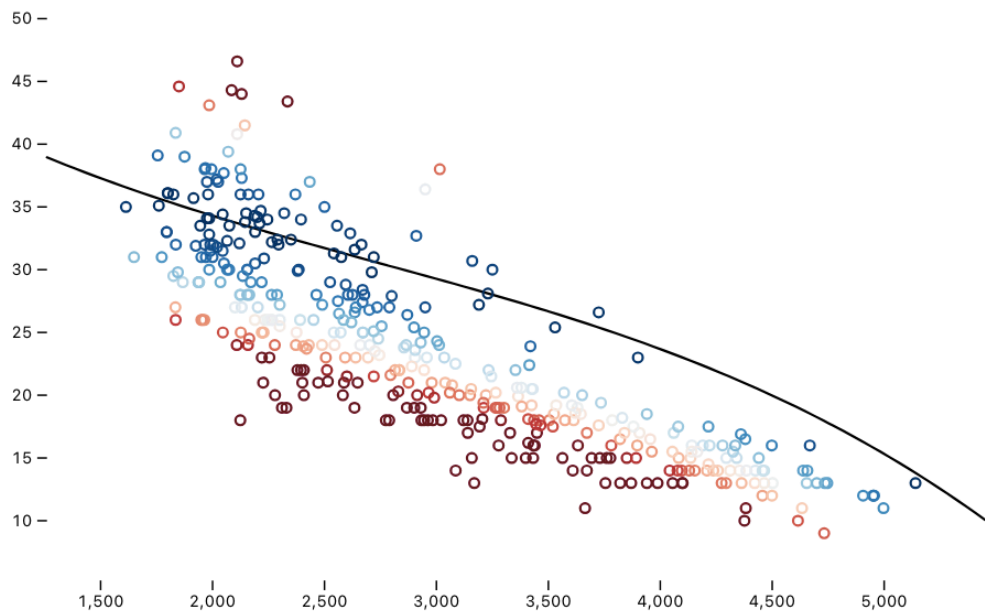


Quadratic decision boundary



Cubic feature transforms

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



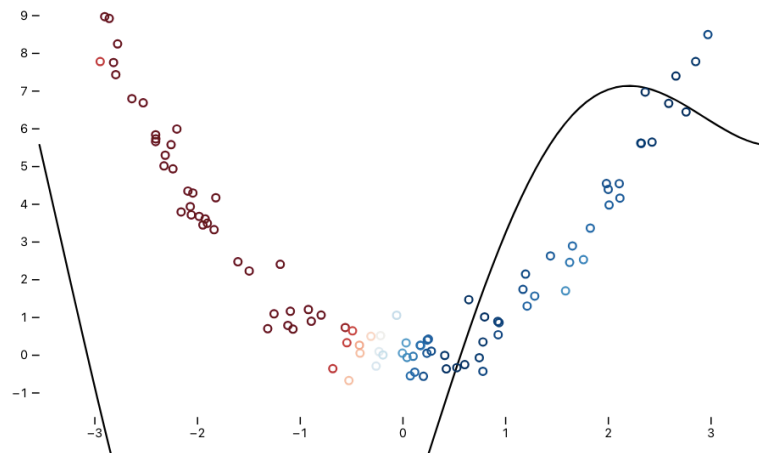
Other feature transforms

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

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

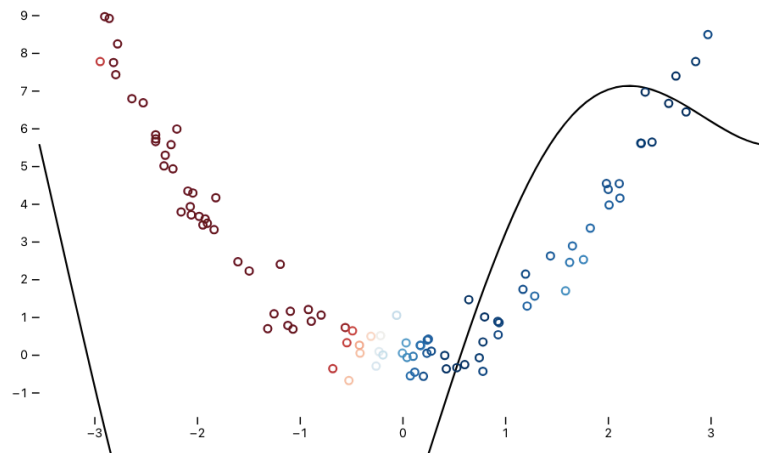
Other feature transforms

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

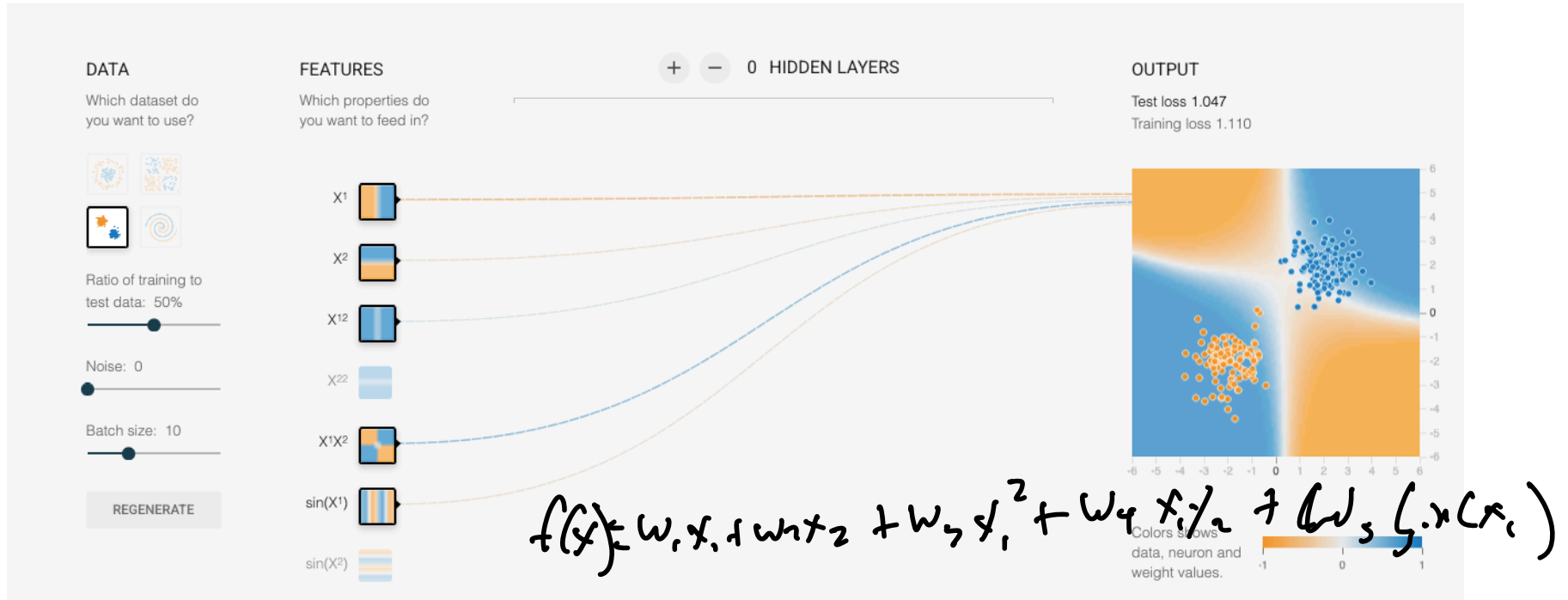


Other feature transforms

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



Back to our viz



How do we choose a transform?

Data = \mathbf{X}, \mathbf{y}



Training data = $\mathbf{X}_{train}, \mathbf{y}_{train}$ Split Test data = $\mathbf{X}_{test}, \mathbf{y}_{test}$



Fit model

$$\mathbf{w} \leftarrow \underset{\mathbf{w}}{\operatorname{argmin}} \operatorname{Loss}(\mathbf{w}, \mathbf{X}_{train}, \mathbf{y}_{train})$$



Evaluate model

$$\operatorname{Loss}(\mathbf{w}, \mathbf{X}_{test}, \mathbf{y}_{test})$$