

# CSE 446

- Linear Regression w/ Basis Functions
- Model Evaluation

---

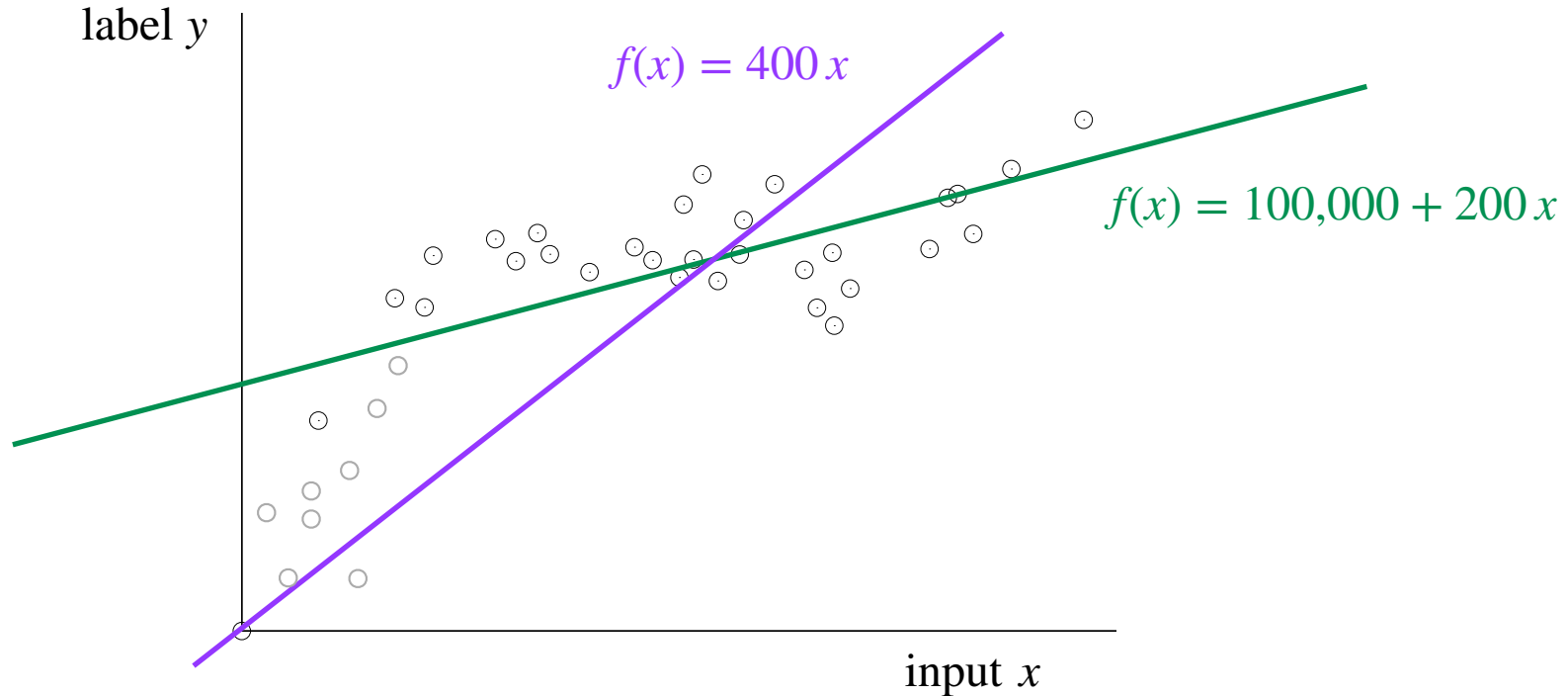
Natasha Jaques

# Polynomial Features

---



# Recap: Linear Regression



- In general high-dimensions, we fit a linear model with intercept  $y_i \simeq w^T x_i + b$ , or equivalently  $y_i = w^T x_i + b + \epsilon_i$  with model parameters  $(w \in \mathbb{R}^d, b \in \mathbb{R})$  that minimizes  $\ell_2$ -loss

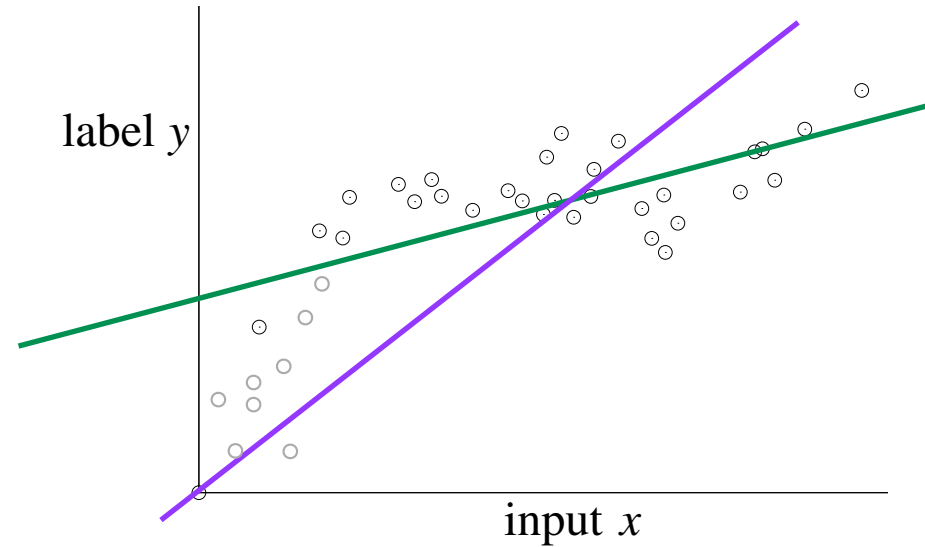
$$\mathcal{L}(w, b) = \sum_{i=1}^n \underbrace{(y_i - (w^T x_i + b))^2}_{\text{error } \epsilon_i} \quad \# \text{ MSE}$$

# Quadratic regression in 1-dimension

- **Data:**  $\mathbf{X} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$ ,  $\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$

- **Linear model with parameter  $(b, w_1)$ :**

- $\hat{y}_i = \underline{b} + \underline{w_1 x_i}$



# Quadratic regression in 1-dimension

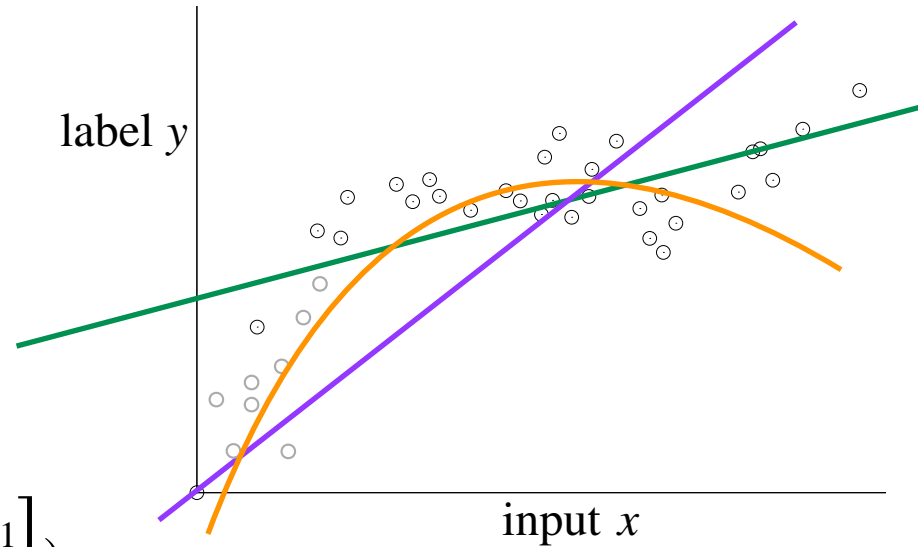
- **Data:**  $\mathbf{X} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$ ,  $\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$

- **Linear model with parameter  $(b, w_1)$ :**

- $\hat{y}_i = b + w_1 x_i$

- **Quadratic model with parameter  $(b, w = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix})$ :**

- $\hat{y}_i = b + w_1 x_i + w_2 x_i^2$



$$h(x_i) = \begin{bmatrix} 1 \\ x_i \\ x_i^2 \end{bmatrix} \quad \hat{y} = h(x_i)^T \begin{bmatrix} b \\ w_1 \\ w_2 \\ \vdots \end{bmatrix}$$

# Still linear regression, but on quadratic features

# Linear layer fit to non-linear features

# Quadratic regression in 1-dimension

- **Data:**  $\mathbf{X} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$ ,  $\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$

- **Linear model with parameter  $(b, w_1)$ :**

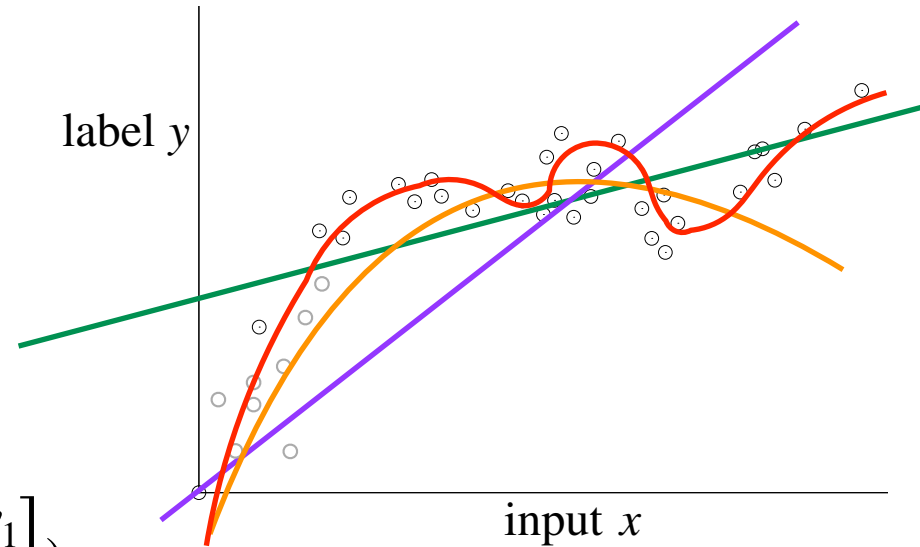
- $\hat{y}_i = b + w_1 x_i$

- **Quadratic model with parameter  $(b, w = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix})$ :**

- $\hat{y}_i = b + w_1 x_i + w_2 x_i^2$

- **Degree-p polynomial model with  $p$  parameters**

- $\hat{y}_i = b + w_1 x_i + w_2 x_i^2 + \dots + w_p x_i^p$



$$h(x_i) = \begin{bmatrix} 1 \\ x_i \\ x_i^2 \end{bmatrix} \quad \hat{y} = h(x_i)^T \begin{bmatrix} b \\ w_1 \\ w_2 \\ \vdots \\ w_p \end{bmatrix}$$

# Still linear regression, but on quadratic features

# Linear layer fit to non-linear features

# Quadratic regression in 1-dimension

- **Data:**  $\mathbf{X} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$ ,  $\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$

- **Linear model with parameter  $(b, w_1)$ :**

- $\hat{y}_i = b + w_1 x_i$

- **Quadratic model with parameter  $(b, w = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix})$ :**

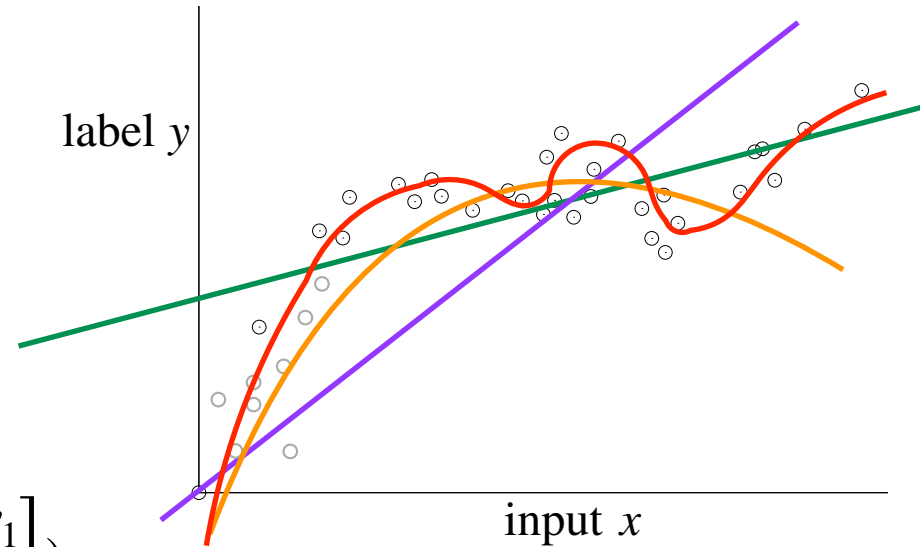
- $\hat{y}_i = b + w_1 x_i + w_2 x_i^2$

- **Degree-p polynomial model with  $p$  parameters**

- $\hat{y}_i = b + w_1 x_i + w_2 x_i^2 + \dots + w_p x_i^p$

- **General p-features with parameter  $w = \begin{bmatrix} w_1 \\ \vdots \\ w_p \end{bmatrix}$ :**

- $\hat{y}_i = \langle w, h(x_i) \rangle$  where  $h : \mathbb{R} \rightarrow \mathbb{R}^p$

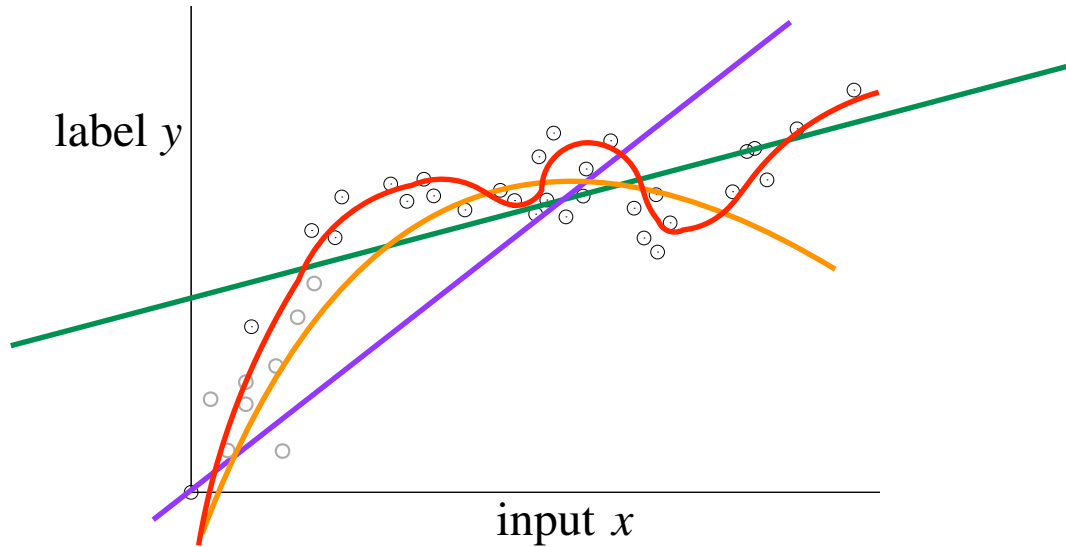


$$h(x_i) = \begin{bmatrix} 1 \\ x_i \\ x_i^2 \end{bmatrix} \quad \hat{y} = h(x_i)^T \begin{bmatrix} b \\ w_1 \\ w_2 \\ \vdots \\ w_p \end{bmatrix}$$

# Still linear regression, but on quadratic features

# Linear layer fit to non-linear features

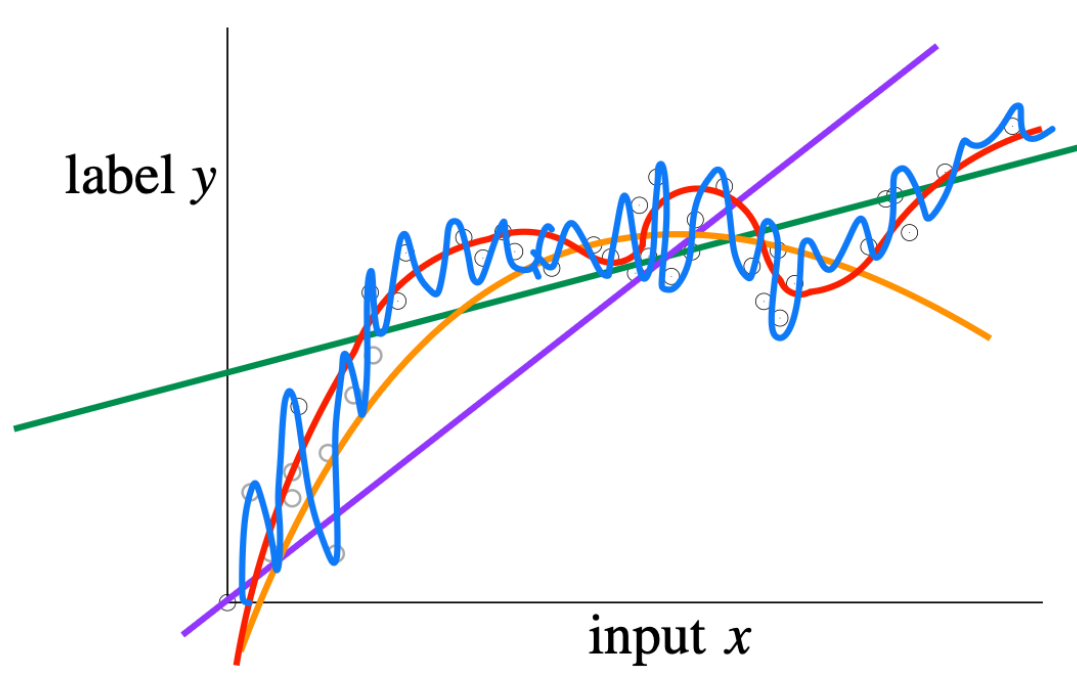
# Quadratic regression in 1-dimension



Which line gives us the best fit?

- **Green & purple?** Can't capture variance. **Underfitting**
- **Orange?** Can't capture variance. **Underfitting**
- **Red?**  
Just right?  
Could we take it too far?

# Quadratic regression in 1-dimension



Which line gives us the best fit?

- **Green** & **purple**? Can't capture variance. **Underfitting**
- **Orange**? Can't capture variance. **Underfitting**
- **Red**? Just right?
- **Blue**? Fits random fluctuations / measurement errors in training dataset. **Overfitting.**

# Quadratic regression in 1-dimension

- **Data:**  $\mathbf{X} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$ ,  $\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$

- **General p-features with parameter**  $w = \begin{bmatrix} w_1 \\ \vdots \\ w_p \end{bmatrix}$  :
  - $\hat{y}_i = \langle w, h(x_i) \rangle$  where  $h : \mathbb{R} \rightarrow \mathbb{R}^p$

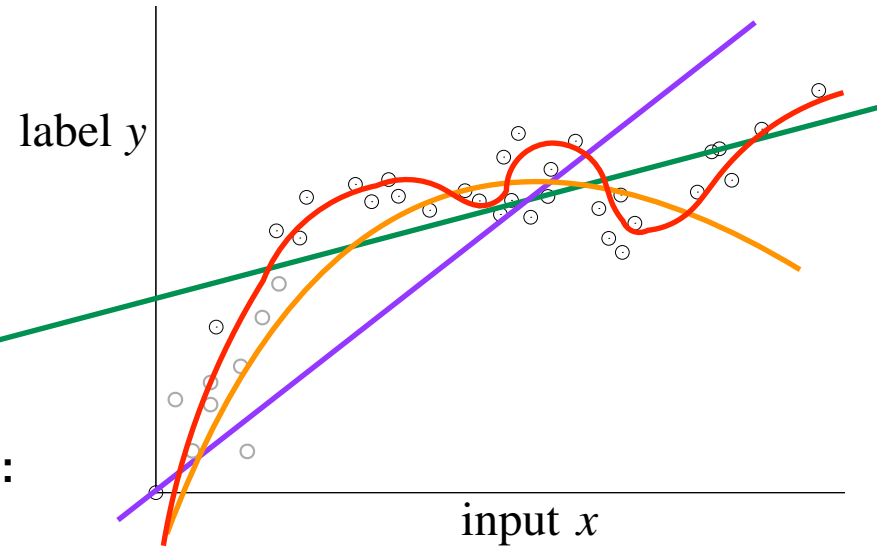
Note:  $h$  can be arbitrary non-linear functions!

$$h(x) = \left[ \log(x), x^2, \sin(x), \sqrt{x} \right]^T$$

# Transform features however

# How will I know if I have good features?

PLOT PLOT PLOT PLOT PLOT



## Breakout discussion:

- When would you use  $\sin(x)$ ?

# Cyclic data e.g. weather

- How should I transform the house feature: zip code?

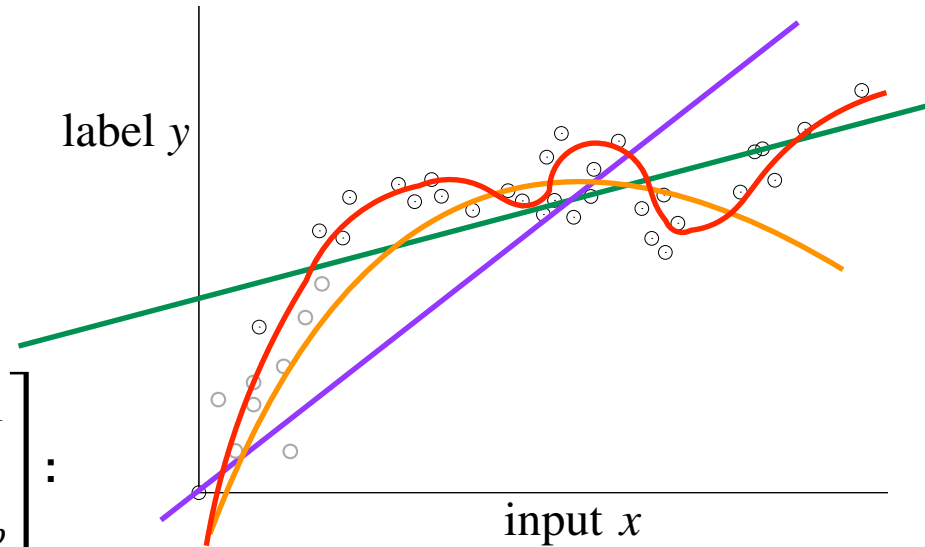
# Lat / long

**This is the art of feature engineering**

# Quadratic regression in 1-dimension

- **Data:**  $\mathbf{X} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$ ,  $\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$

- **General p-features with parameter**  $w = \begin{bmatrix} w_1 \\ \vdots \\ w_p \end{bmatrix}$  :
  - $\hat{y}_i = \langle w, h(x_i) \rangle$  where  $h : \mathbb{R} \rightarrow \mathbb{R}^p$



How do we learn  $w$ ?

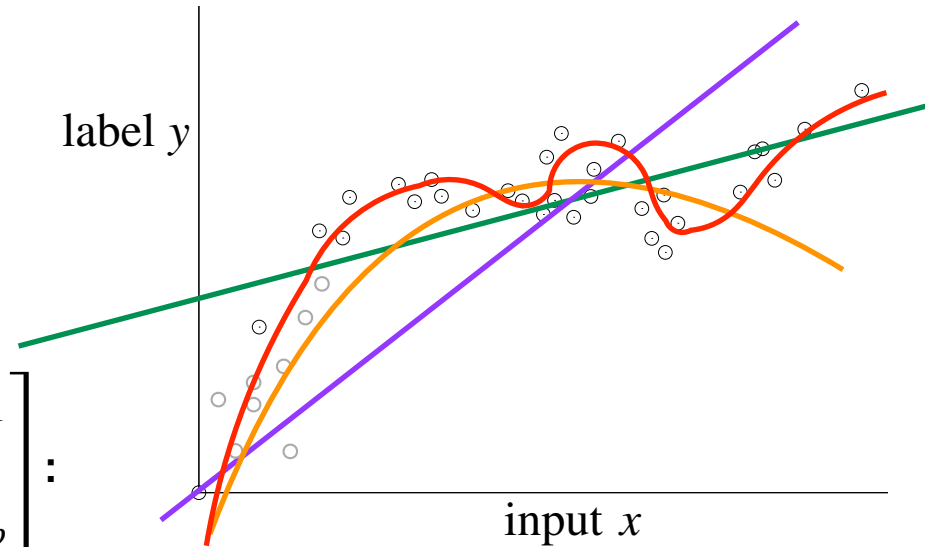
# Quadratic regression in 1-dimension

- **Data:**  $\mathbf{X} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$ ,  $\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$

- **General p-features with parameter  $w =$**

$$\begin{bmatrix} w_1 \\ \vdots \\ w_p \end{bmatrix} :$$

- $\hat{y}_i = \langle w, h(x_i) \rangle$  where  $h : \mathbb{R} \rightarrow \mathbb{R}^p$



How do we learn  $w$ ?

$$\mathbf{H} = \begin{bmatrix} - - h(x_1)^T - - \\ \vdots \\ - - h(x_n)^T - - \end{bmatrix} \in \mathbb{R}^{n \times p}$$

$$\hat{w} = \arg \min_w \|\mathbf{H}w - \mathbf{y}\|_2^2$$

For a new test point  $x$ , predict  
 $\hat{y} = \langle \hat{w}, h(x) \rangle$

$$\hat{w}_{\text{MLE}} = (\mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T \mathbf{y}$$