

Regression

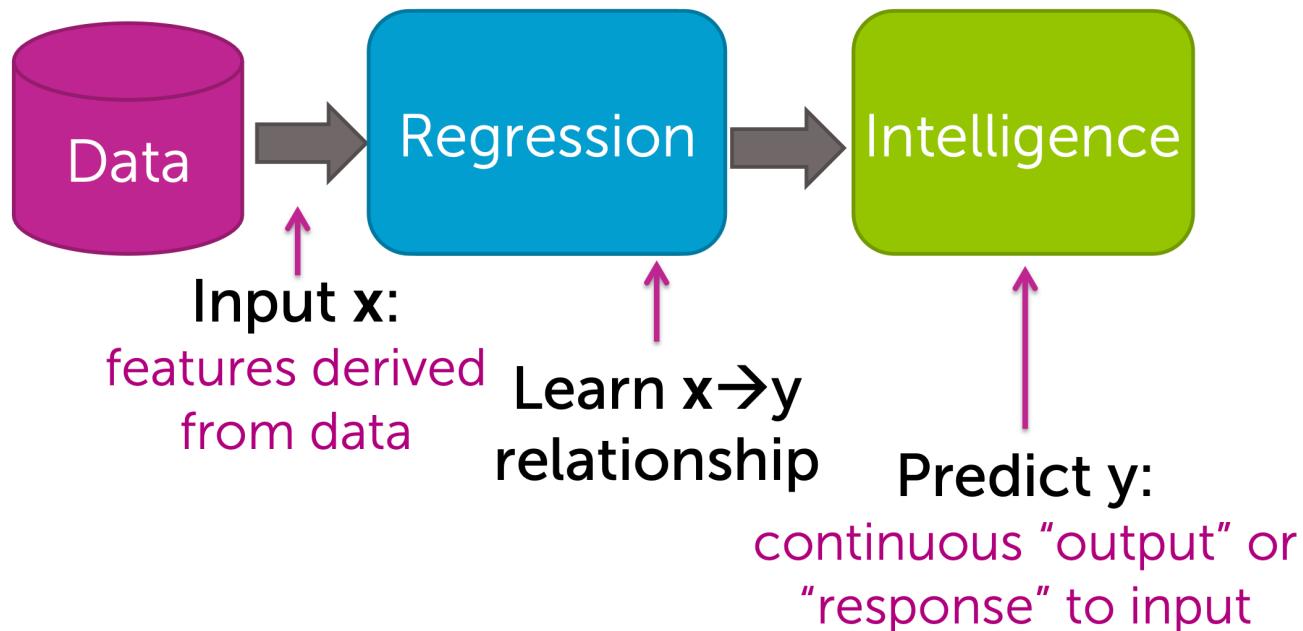
K Kotecha

<https://www.linkedin.com/in/ketankotecha/>

- Material taken in this slides are from various sources. Its only used for education purpose.

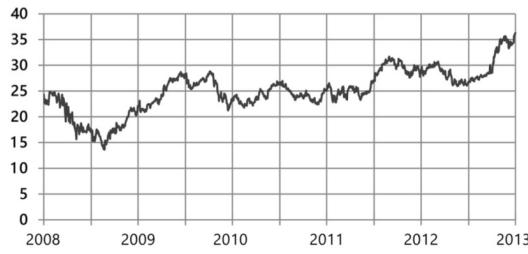
What is regression?

From features to predictions



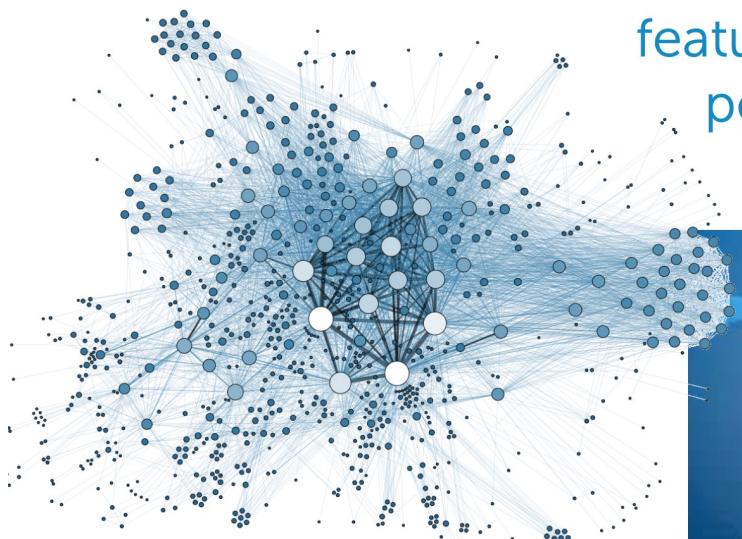
Stock prediction

- Predict the price of a stock (y)
- Depends on $x =$
 - Recent history of stock price
 - News events
 - Related commodities

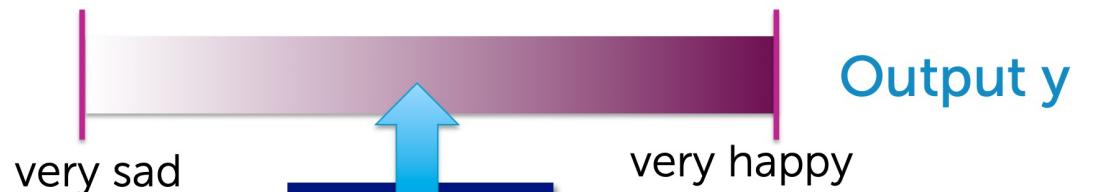


Tweet popularity

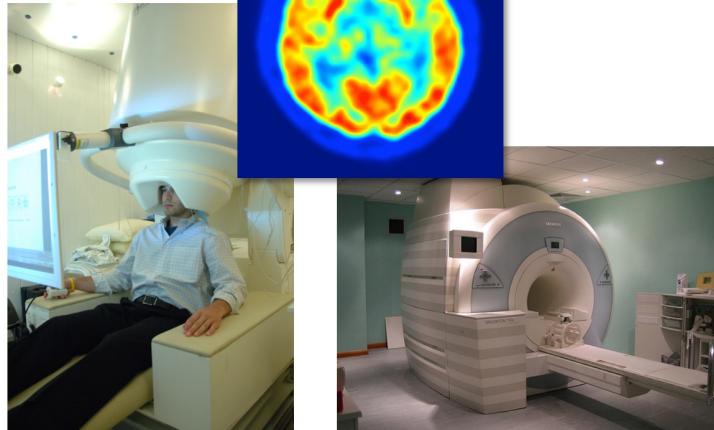
- How many people will retweet your tweet? (y)
- Depends on $x = \# \text{ followers}$,
 $\# \text{ of followers of followers}$,
 $\text{features of text tweeted}$,
 $\text{popularity of hashtag}$,
 $\# \text{ of past retweets}, \dots$



Reading your mind

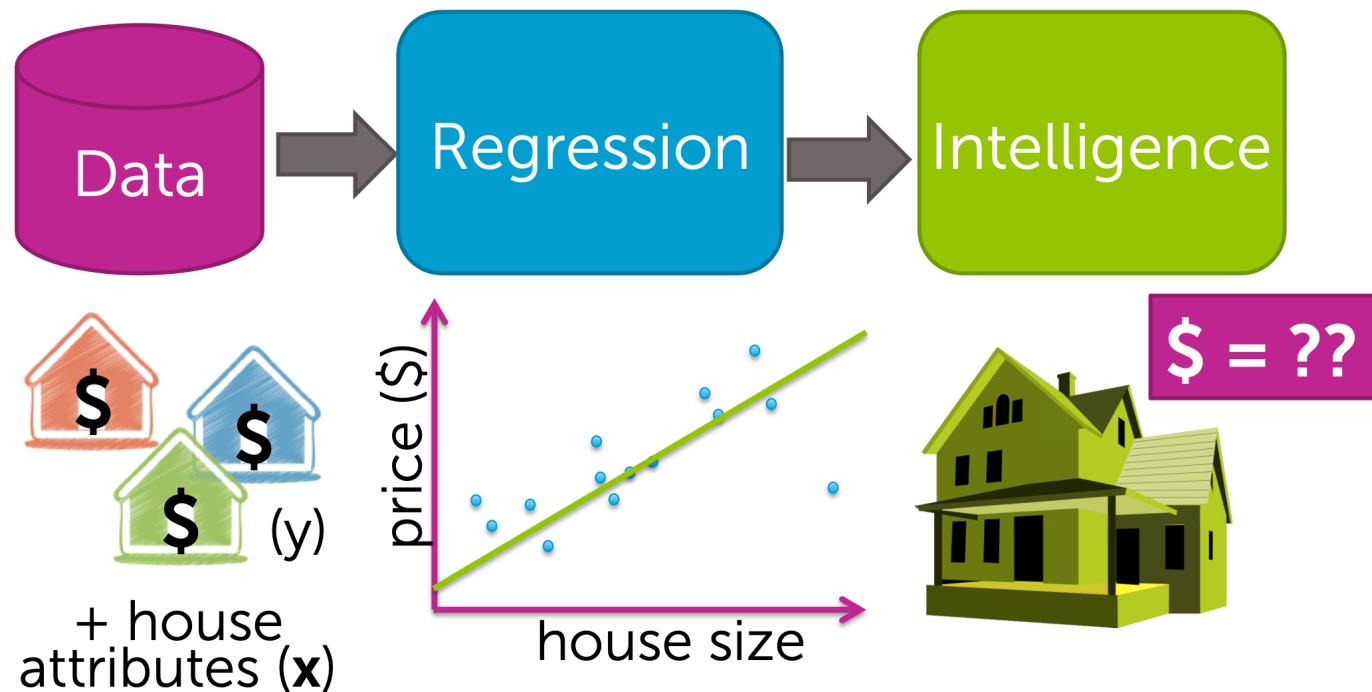


Output y



Inputs x are
brain region
intensities

Case Study: Predicting house prices



Simple Linear Regression

Constant

Coefficient

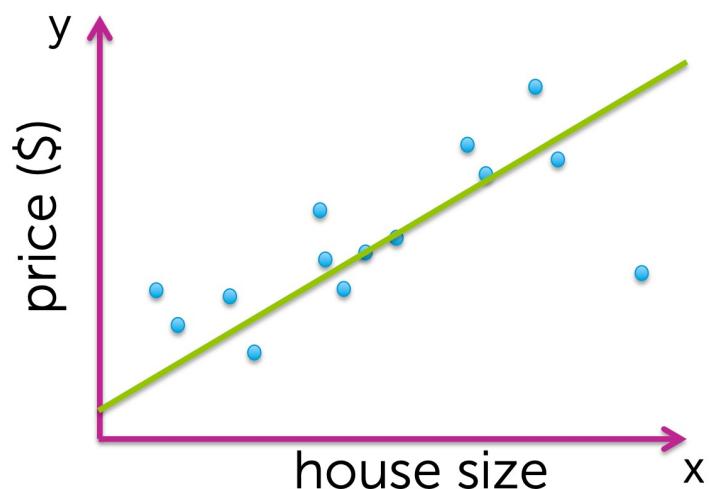
$$y = b_0 + b_1 * x_1$$

Dependent variable (DV)

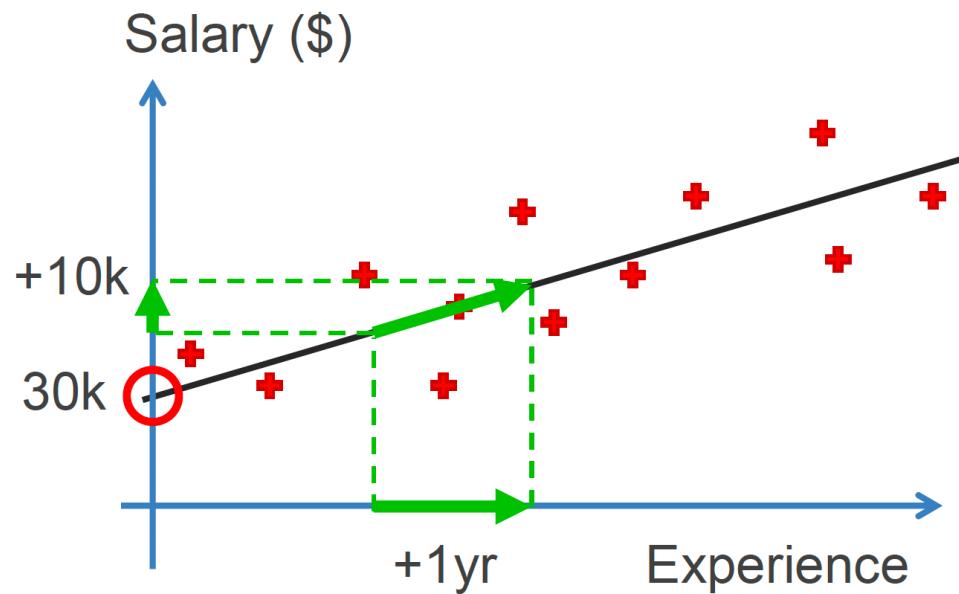
Independent variable (IV)

Simple Regression

What makes it simple?
1 input and just fit a line to data



Simple Linear Regression:



$$y = b_0 + b_1 * x$$

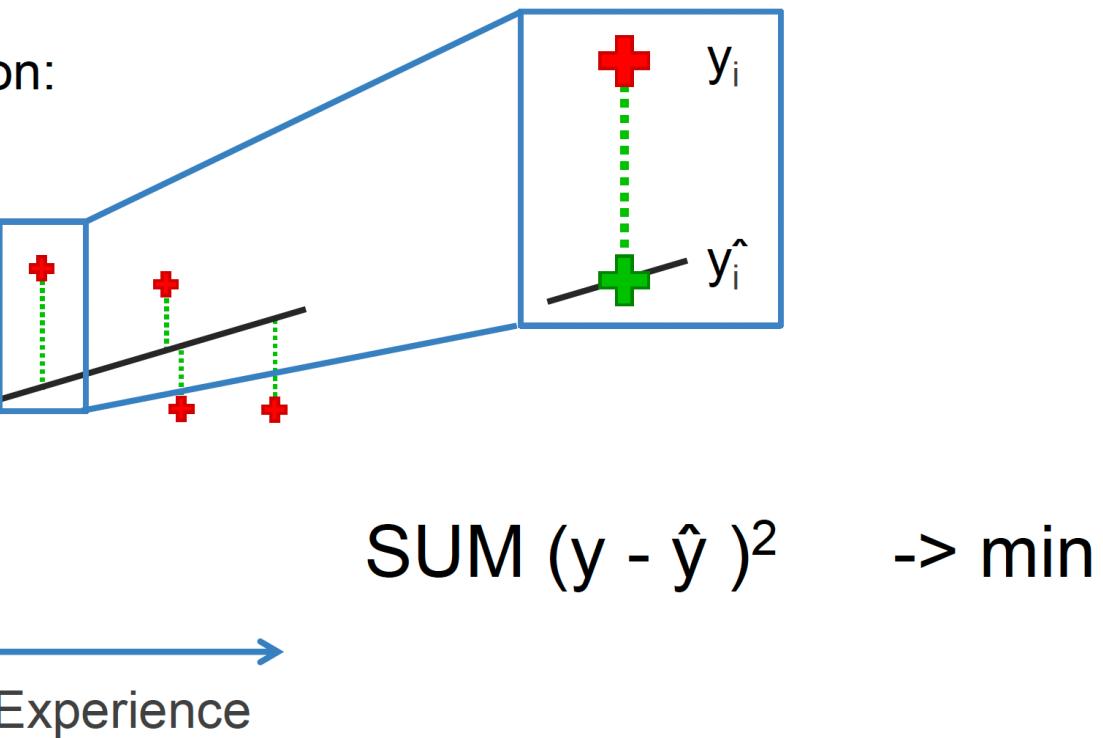
↓
Salary = b₀ + b₁ * Experience

Simple Linear Regression:

Salary (\$)



Experience



How much is my house worth?



Look at recent sales in my neighborhood

- How much did they sell for?



Data



input *output*
 $(x_1 = \text{sq.ft.}, y_1 = \$)$



$(x_2 = \text{sq.ft.}, y_2 = \$)$



$(x_3 = \text{sq.ft.}, y_3 = \$)$



$(x_4 = \text{sq.ft.}, y_4 = \$)$



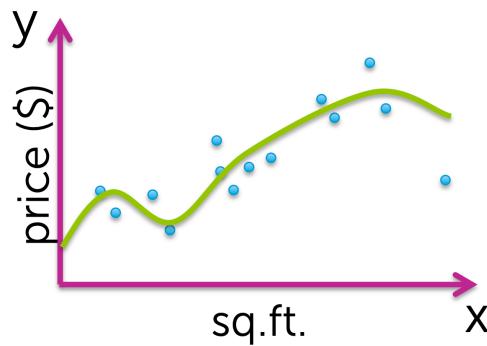
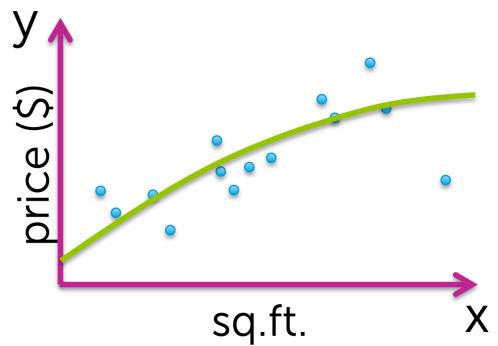
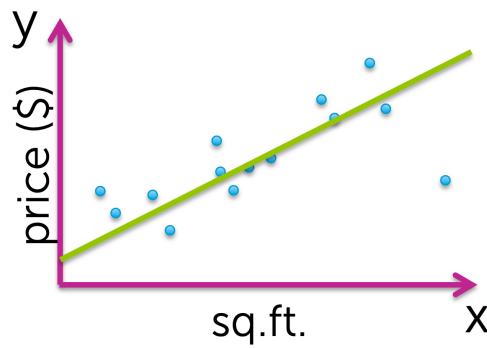
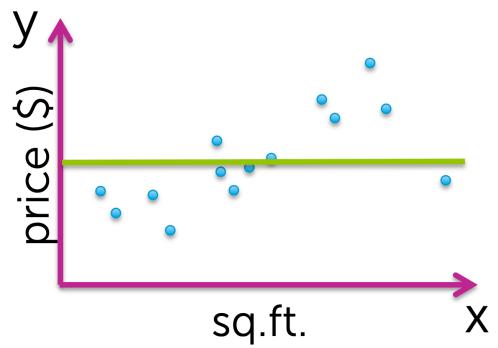
$(x_5 = \text{sq.ft.}, y_5 = \$)$

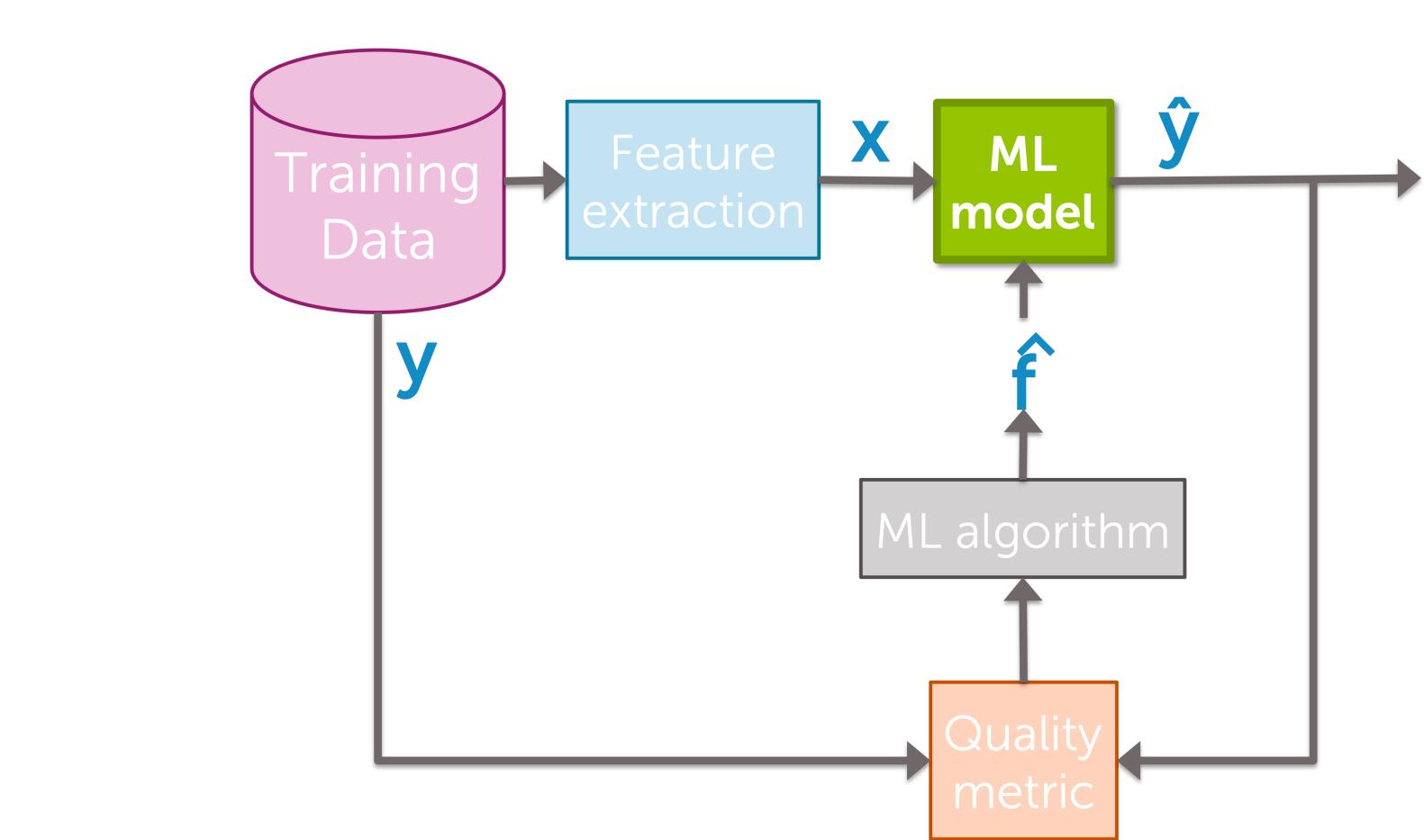
.

Input vs. Output:

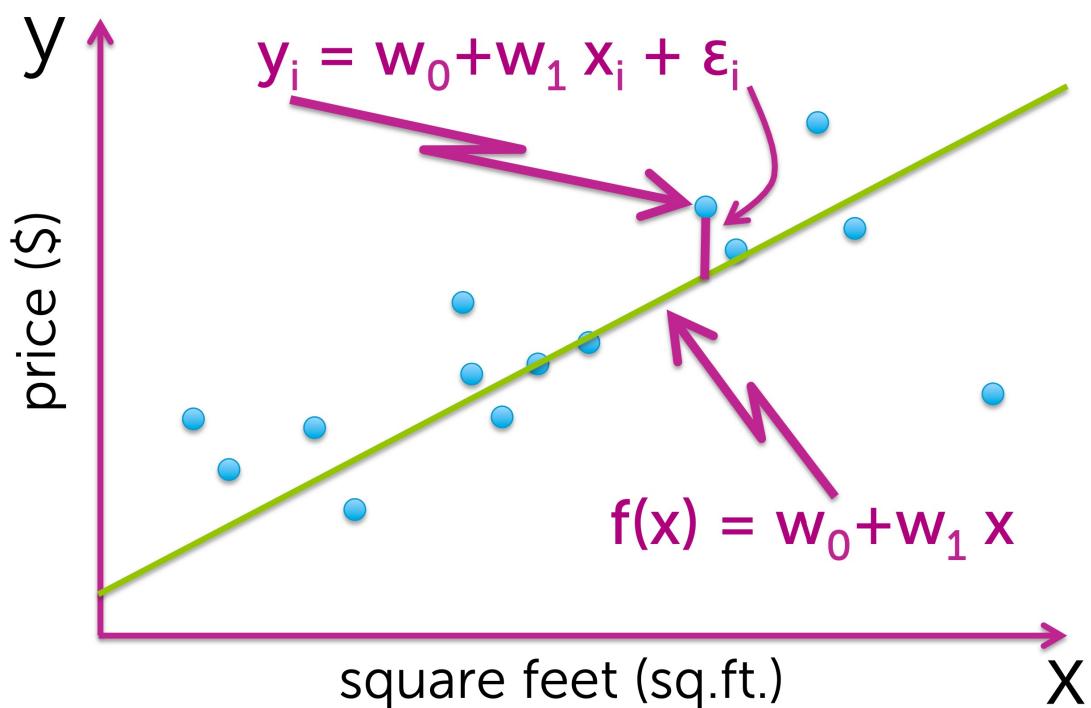
- y is the quantity of interest
- assume y can be predicted from x

Task 1– Which model $f(x)$?

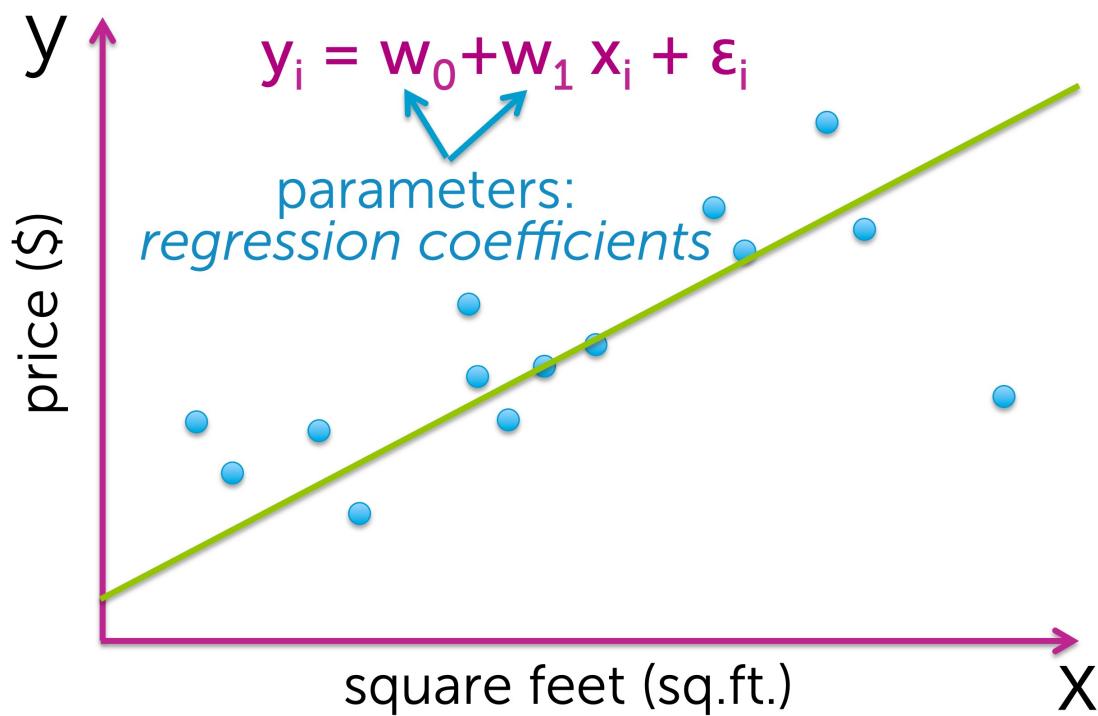


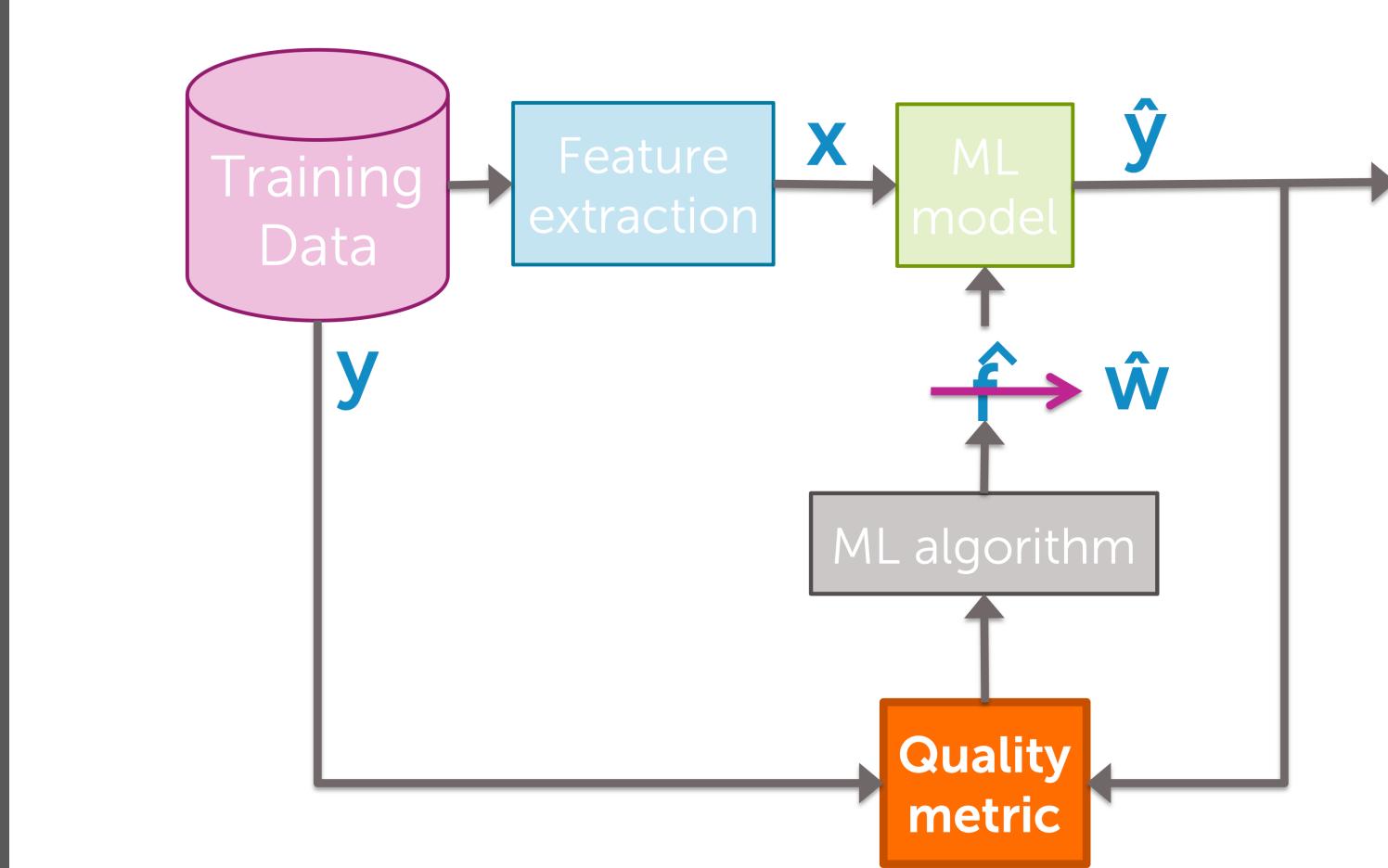


Simple linear regression model



Simple linear regression model





Living area (feet ²)	Price (1000\$s)
2104	400
1600	330
2400	369
1416	232
3000	540
:	:

Simple
Linear
Regression

$$y = b_0 + b_1 * x_1$$

Multiple Linear Regression

Living area (feet ²)	#bedrooms	Price (1000\$s)
2104	3	400
1600	3	330
2400	3	369
1416	2	232
3000	4	540
:	:	:

Multiple
Linear
Regression

$$y = b_0 + b_1 * x_1 + b_2 * x_2$$

Multivariate Regression

Size (feet ²)	Number of bedrooms	Number of floors	Age of home (years)	Price (\$1000)
2104	5	1	45	460
1416	3	2	40	232
1534	3	2	30	315
852	2	1	36	178

Multiple
Linear
Regression

$$y = b_0 + b_1 * x_1 + b_2 * x_2 + \dots + b_n * x_n$$

Simple Linear
Regression

$$y = b_0 + b_1 * x_1$$

Multiple Linear
Regression

$$y = b_0 + b_1 * x_1 + b_2 * x_2 + \dots + b_n * x_n$$

Dependent variable (DV)

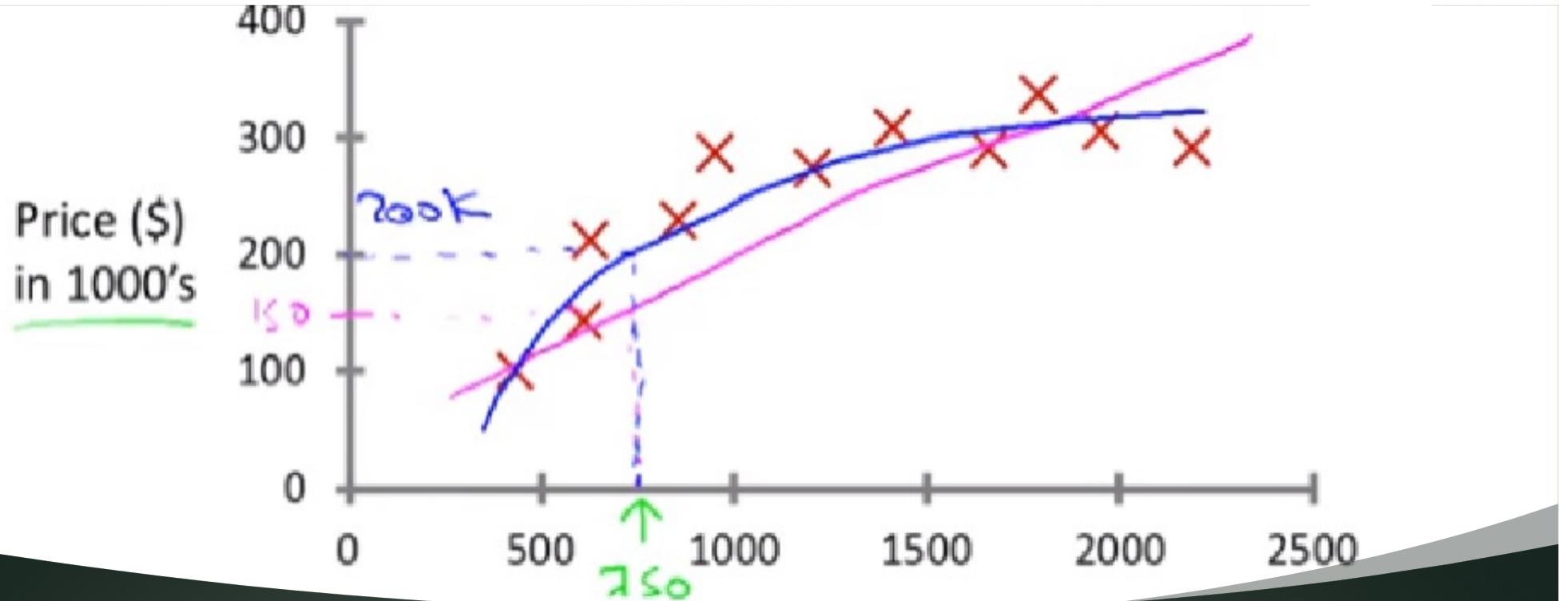
Independent variables (IVs)

Constant

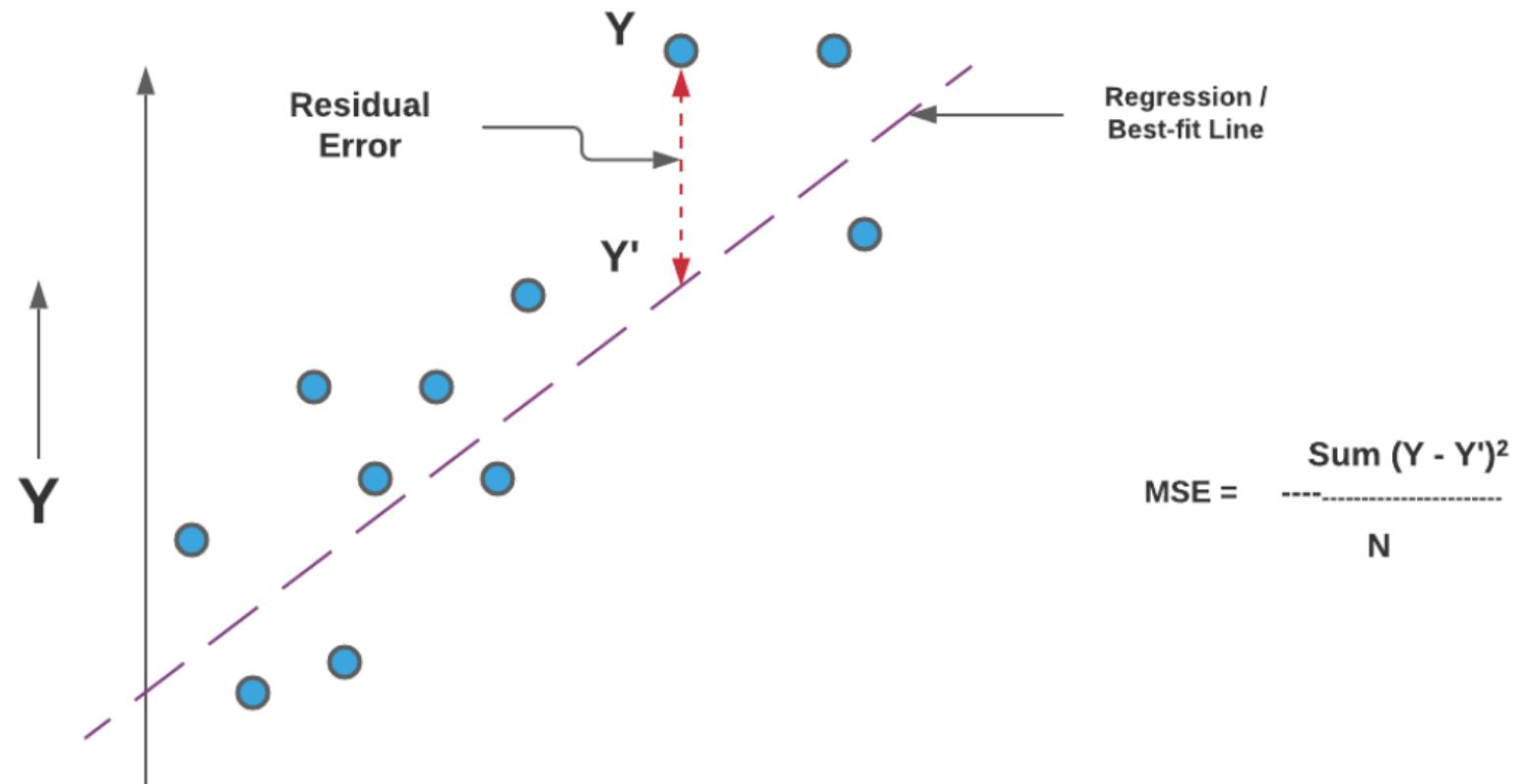
Coefficients

```
graph TD; DV[Dependent variable (DV)] --> y1[y = b0 + b1*x1]; IVs[Independent variables (IVs)] --> y2[y = b0 + b1*x1 + b2*x2 + ... + bn*xn]; Constant[Constant] --> b0[b0]; Coefficients[Coefficients] --> b1[b1*x1]; Coefficients --> b2[b2*x2]; Coefficients --> bn[bn*xn]
```

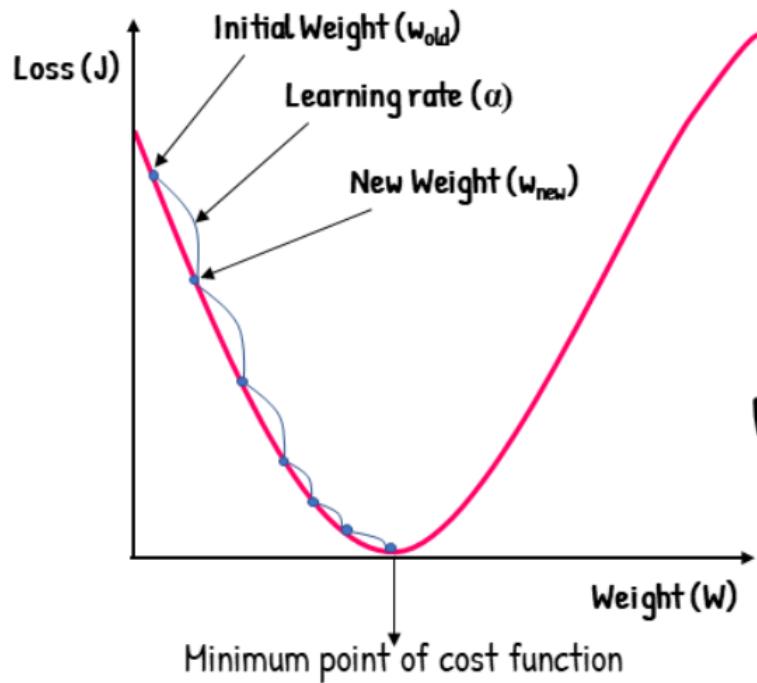
Gradient Descent Algorithm



Housing Price Prediction (Recap)



Gradient Descent



$$w_{new} = w_{old} - \alpha \frac{\delta J}{\delta w}$$

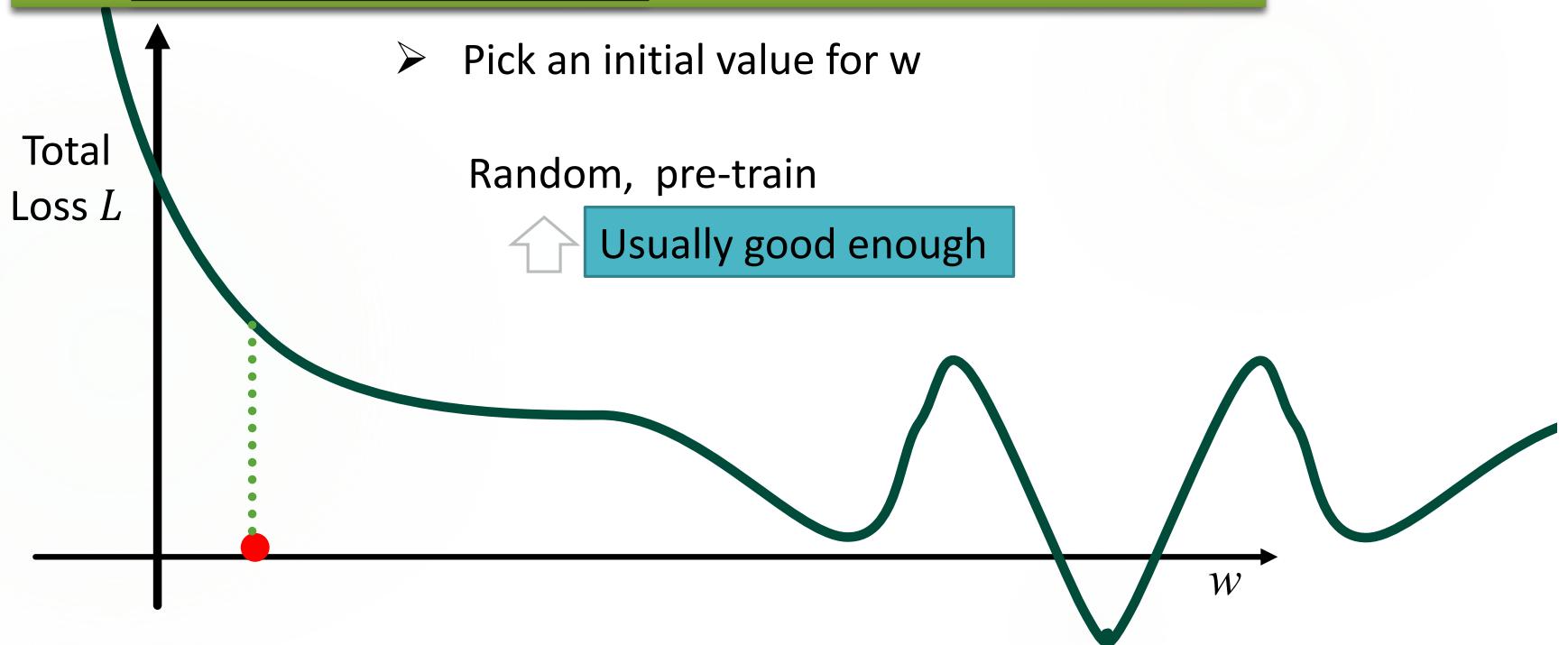
Regression loss functions

Mean Bias Error	Captures average bias in prediction. But is rarely used for training.	$\mathcal{L}_{MBE} = \frac{1}{N} \sum_{i=1}^N (y_i - f(x_i))$
Mean Absolute Error	Measures absolute average bias in prediction. Also called L1 Loss.	$\mathcal{L}_{MAE} = \frac{1}{N} \sum_{i=1}^N y_i - f(x_i) $
Mean Squared Error	Average squared distance between actual and predicted. Also called L2 Loss.	$\mathcal{L}_{MSE} = \frac{1}{N} \sum_{i=1}^N (y_i - f(x_i))^2$
Root Mean Squared Error	Square root of MSE. Loss and dependent variable have same units.	$\mathcal{L}_{RMSE} = \sqrt{\frac{1}{N} \sum_{i=1}^N (y_i - f(x_i))^2}$
Huber Loss	A combination of MSE and MAE. It is parametric loss function.	$\mathcal{L}_{\text{Huberloss}} = \begin{cases} \frac{1}{2}(y_i - f(x_i))^2 & : y_i - f(x_i) \leq \delta \\ \delta(y_i - f(x_i) - \frac{1}{2}\delta) & : \text{otherwise} \end{cases}$
Log Cosh Loss	Similar to Huber Loss + non-parametric. But computationally expensive.	$\mathcal{L}_{\text{LogCosh}} = \frac{1}{N} \sum_{i=1}^N \log(\cosh(f(x_i) - y_i))$

Gradient Descent

Network parameters $\theta = \{w_1, w_2, \dots, b_1, b_2, \dots\}$

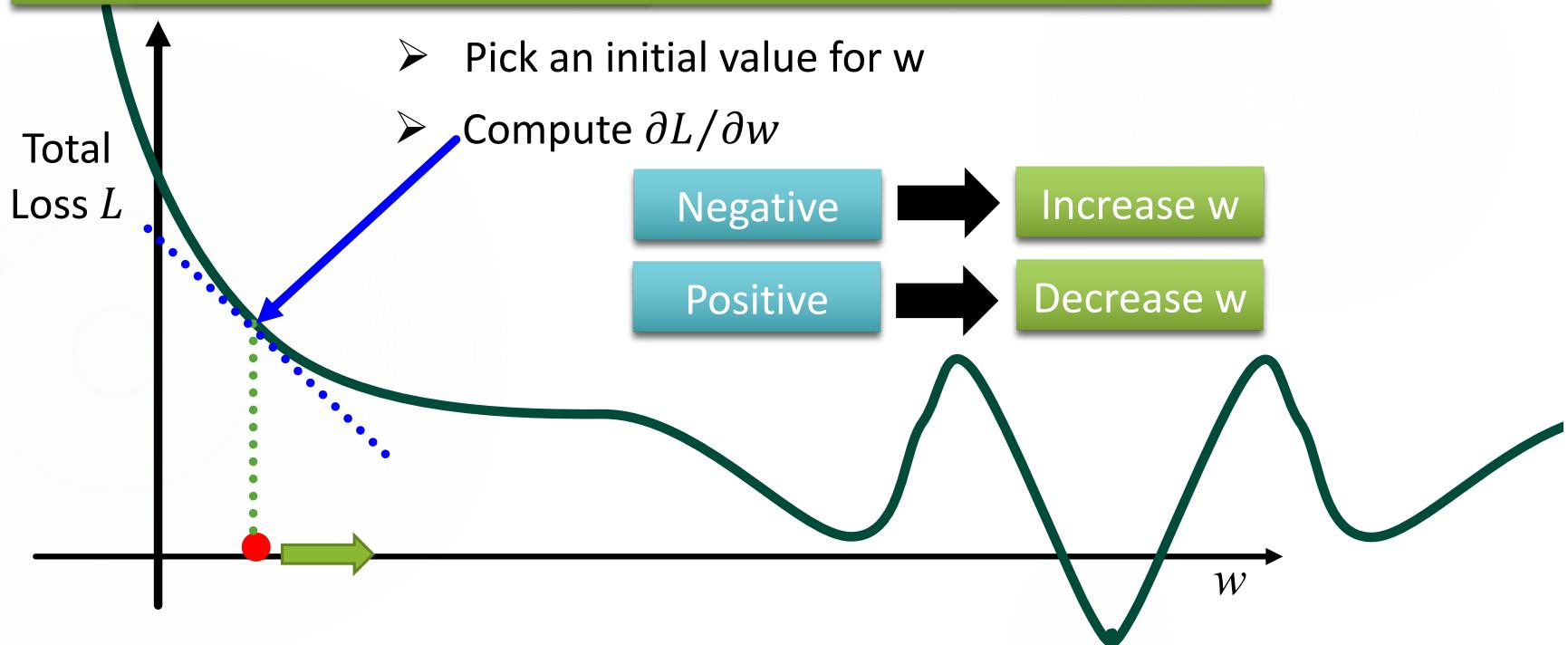
Find network parameters θ^* that minimize total loss L



Gradient Descent

Network parameters $\theta = \{w_1, w_2, \dots, b_1, b_2, \dots\}$

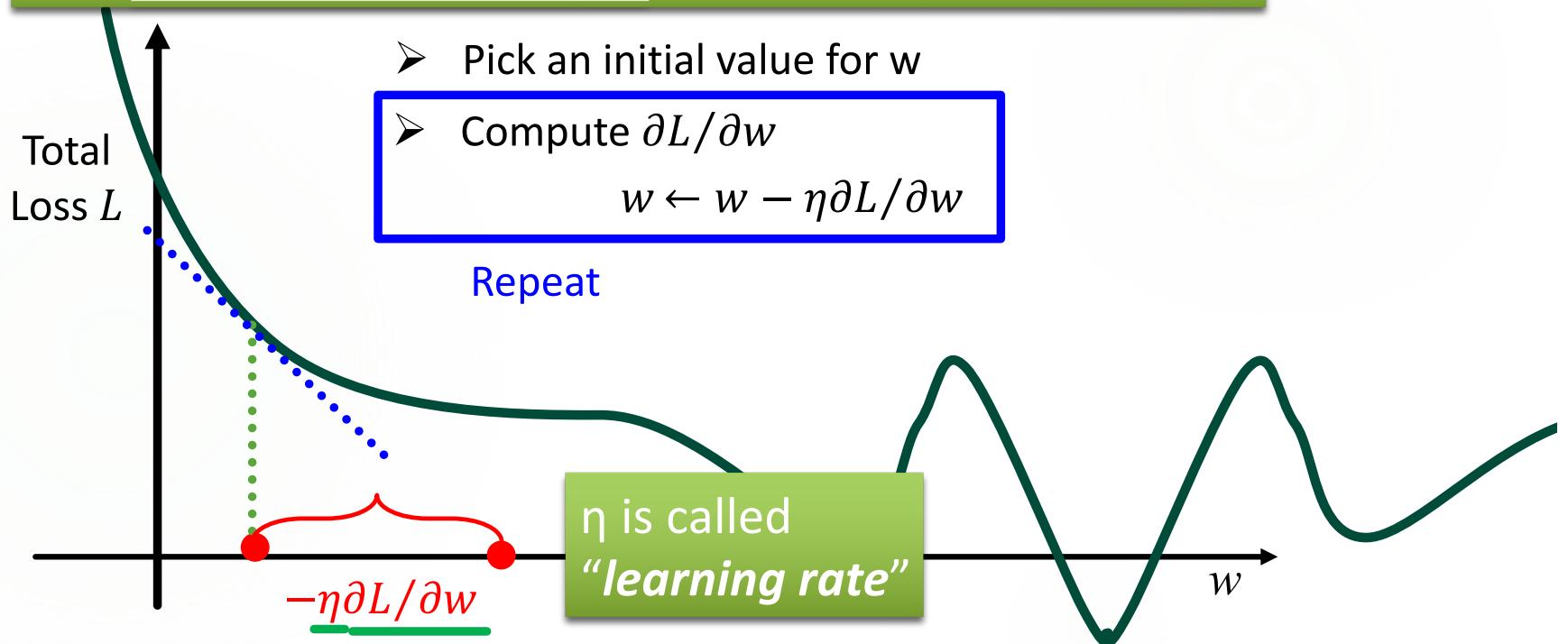
Find network parameters θ^* that minimize total loss L



Gradient Descent

Network parameters $\theta = \{w_1, w_2, \dots, b_1, b_2, \dots\}$

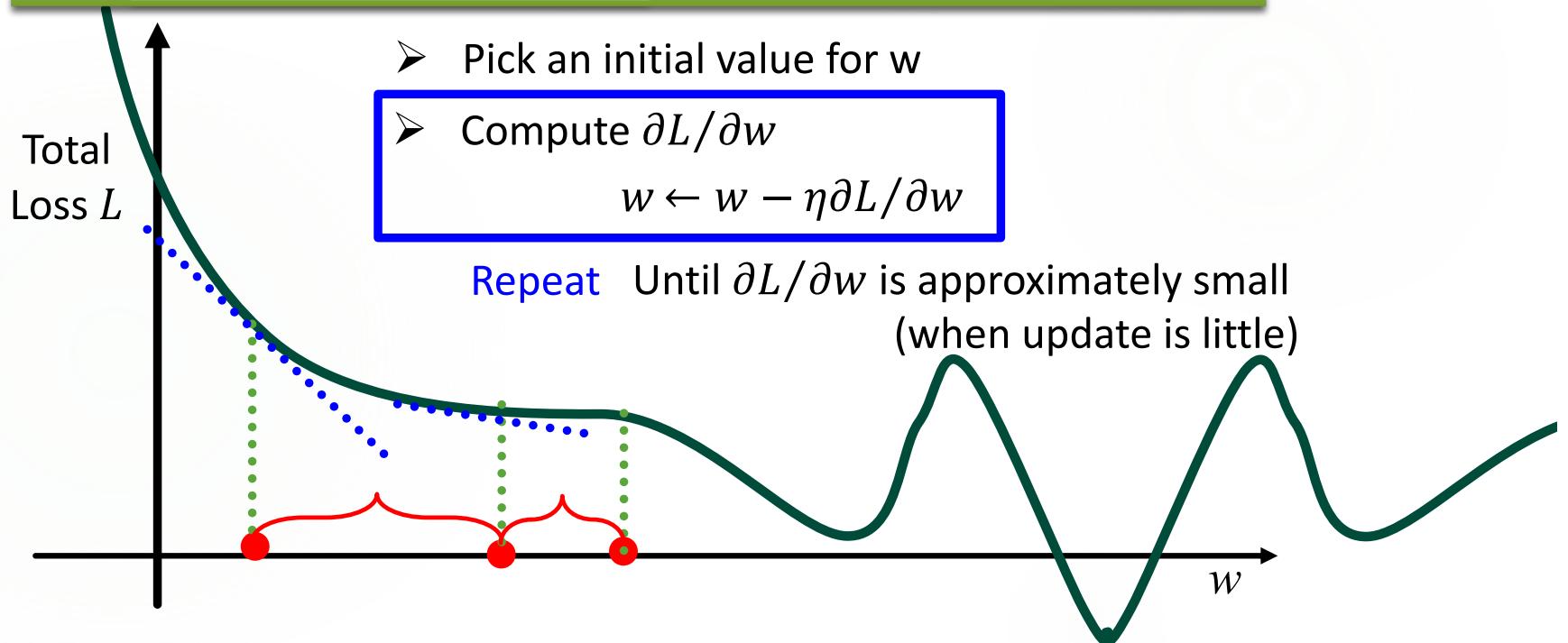
Find network parameters θ^* that minimize total loss L

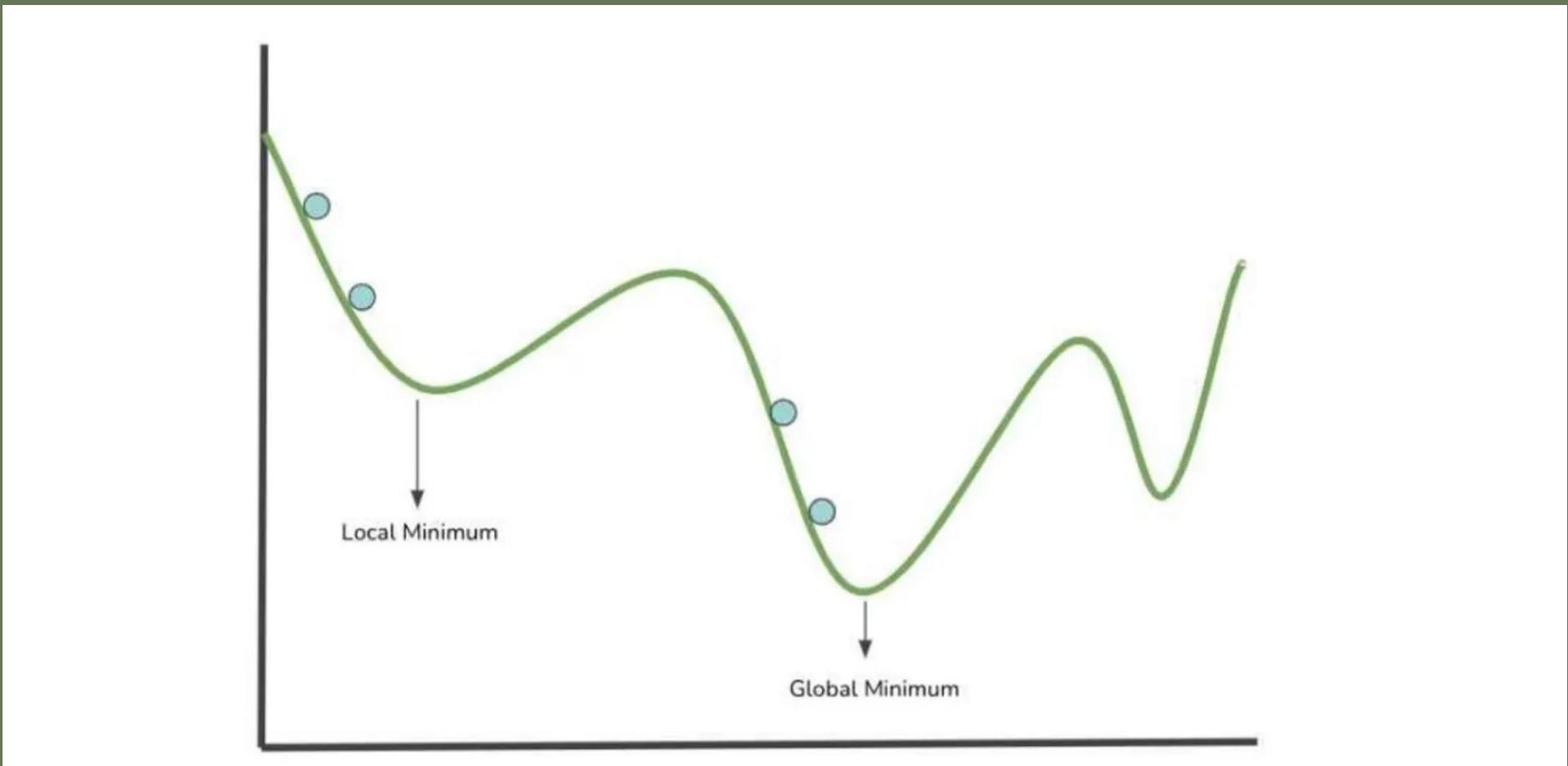


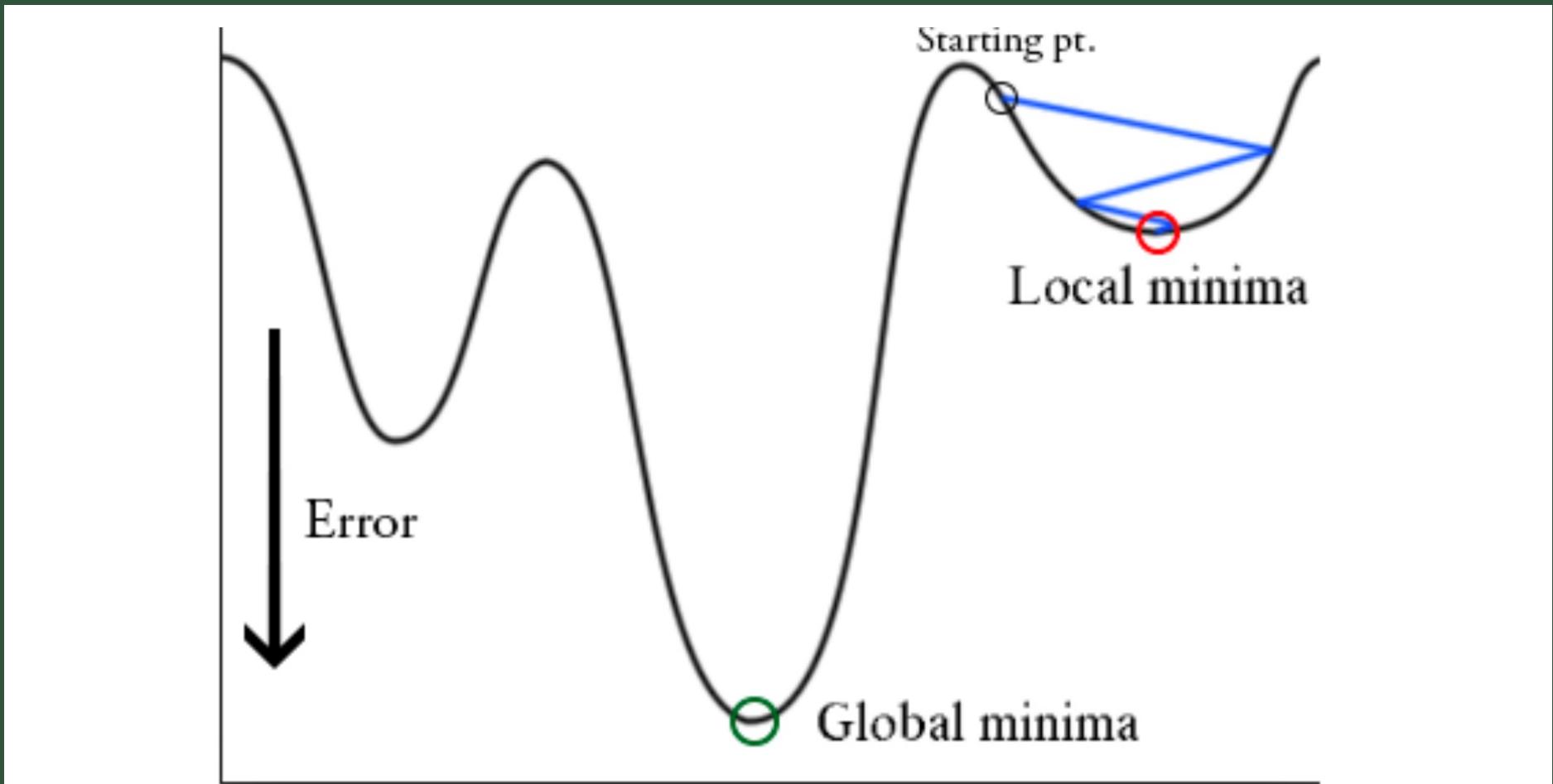
Gradient Descent

Network parameters $\theta = \{w_1, w_2, \dots, b_1, b_2, \dots\}$

Find network parameters θ^* that minimize total loss L







SGD-illustrated

K Kotecha

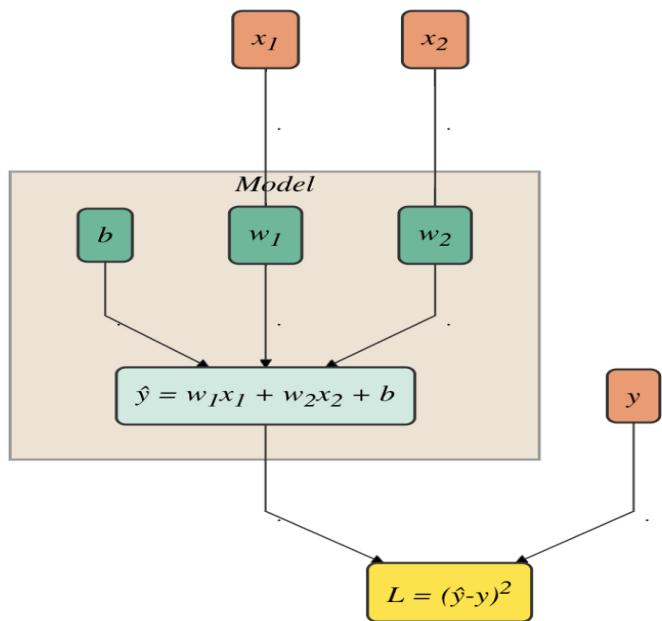
Problem

- | | x1 | x2 | y |
|----|----|----|-----|
| 1) | 4 | 1 | 2 |
| 2) | 2 | 8 | -14 |
| 3) | 1 | 0 | 1 |
| 4) | 3 | 2 | -1 |
| 5) | 1 | 4 | -7 |
| 6) | 6 | 7 | -8 |

$$\hat{y} = w_1x_1 + w_2x_2 + b$$

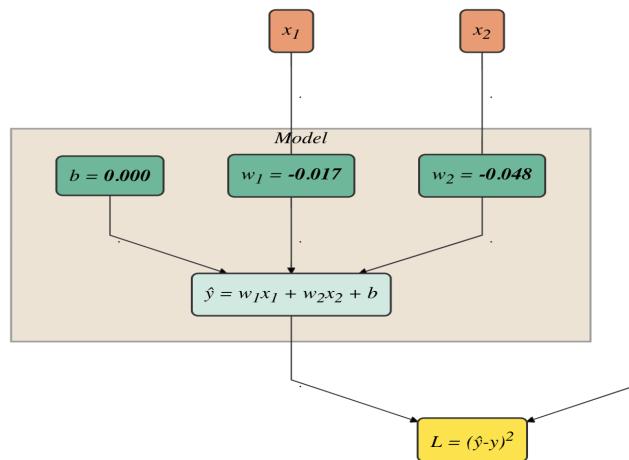
$$\hat{y} = 0.43x_1 - 0.21x_2 + 0.77$$

Representation



Epoch 1 with batch size 1

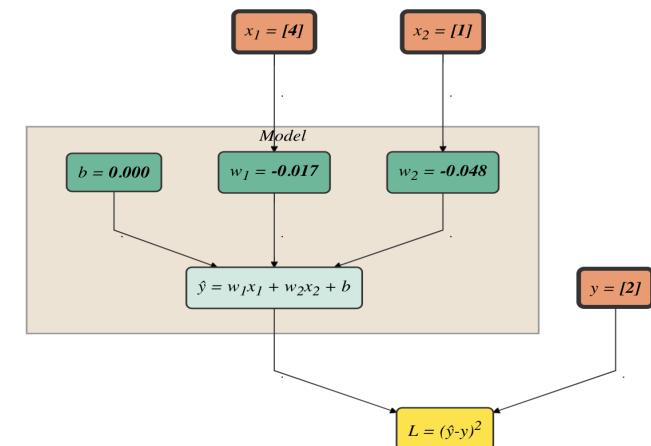
Epoch: - Batch: -



$$w_1 = -0.017$$
$$w_2 = -0.048$$
$$b = 0$$

Y=2

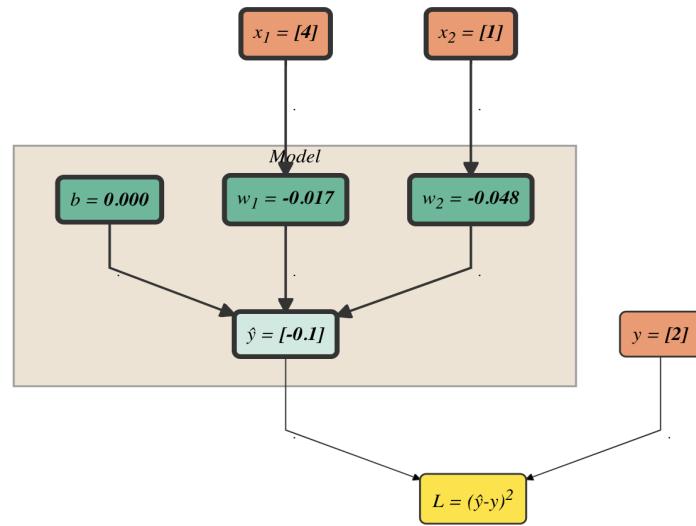
Epoch: 1/1 Batch: 1/6
Forward propagation



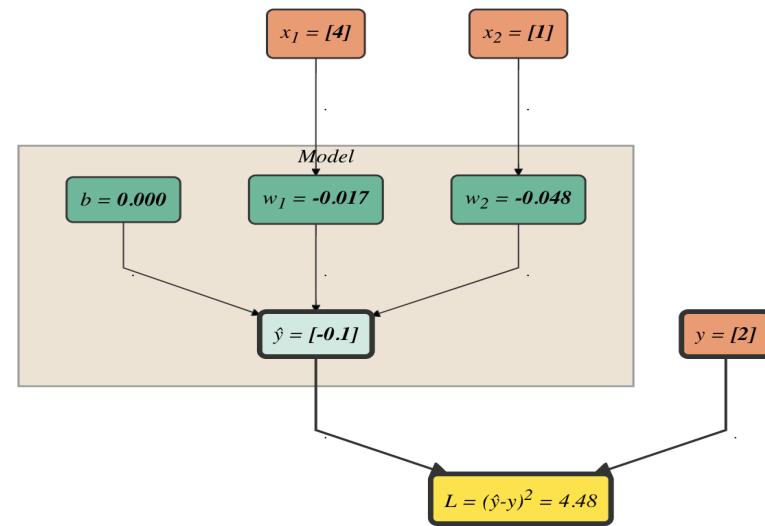
Forward propagation

$$\begin{aligned}\hat{y} &= w_1x_1 + w_2x_2 + b \\ &= (-0.017) \cdot 4 + (-0.048) \cdot 1 + 0.000 \\ &= -0.116\end{aligned}$$

:poch: 1/1 Batch: 1/6
Forward propagation

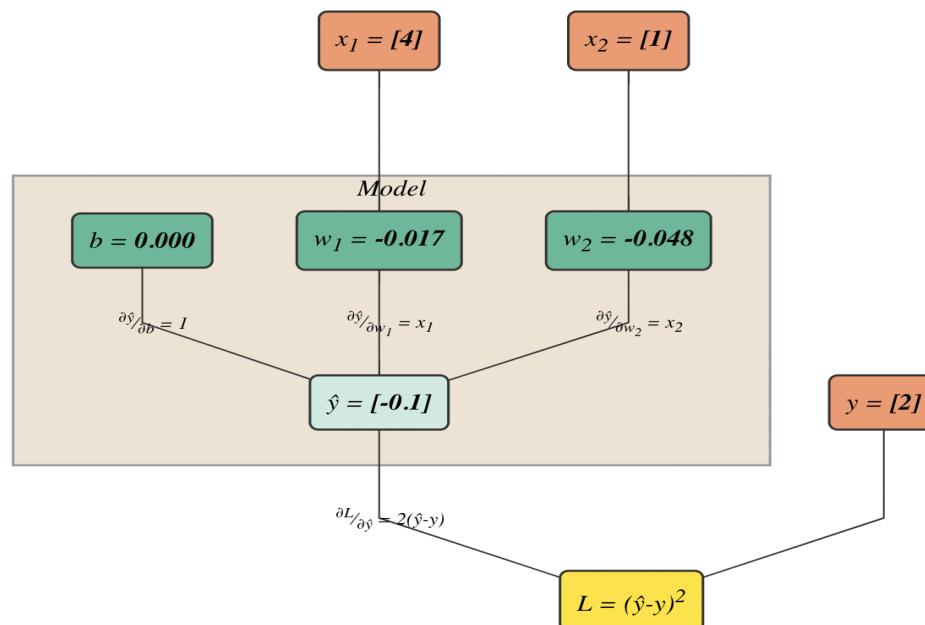


:poch: 1/1 Batch: 1/6
Forward propagation



Derivatives

Epoch: 1/1 Batch: 1/6
Backpropagation



Updating bias

$$\hat{y} = w_1 x_1 + w_2 x_2 + b$$

$$L = (\hat{y} - y)^2$$

$$b' = b - \eta \frac{\partial L}{\partial b}$$

$$\frac{\partial L}{\partial b} = \frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial b}$$

$$\frac{\partial L}{\partial \hat{y}} = 2(\hat{y} - y) \quad \frac{\partial \hat{y}}{\partial b} = 1$$

$$b' = b - \eta \frac{\partial L}{\partial b}$$

$$= b - \eta \left(\frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial b} \right)$$

$$= b - \eta [2(\hat{y} - y) \cdot 1]$$

$$= 0.000 - 0.05[2(-0.116 - 2) \cdot 1]$$

$$= 0.212$$

$$\begin{aligned}x_1 &= 4 \\x_2 &= 1 \\y &= 2 \\\hat{y} &= -0.116\end{aligned}$$

Updating weights

$$\hat{y} = w_1 x_1 + w_2 x_2 + b$$

$$L = (\hat{y} - y)^2$$

$$\begin{aligned}w'_2 &= w_2 - \eta \frac{\partial L}{\partial w_2} \\&= w_2 - \eta \left(\frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial w_2} \right) \\&= w_2 - \eta [2(\hat{y} - y) \cdot x_2] \\&= -0.048 - 0.05[2(-0.116 - 2) \cdot 1] \\&= 0.164\end{aligned}$$

$$\frac{\partial L}{\partial \hat{y}} = 2(\hat{y} - y)$$

$$\frac{\partial \hat{y}}{\partial w_2} = x_2$$

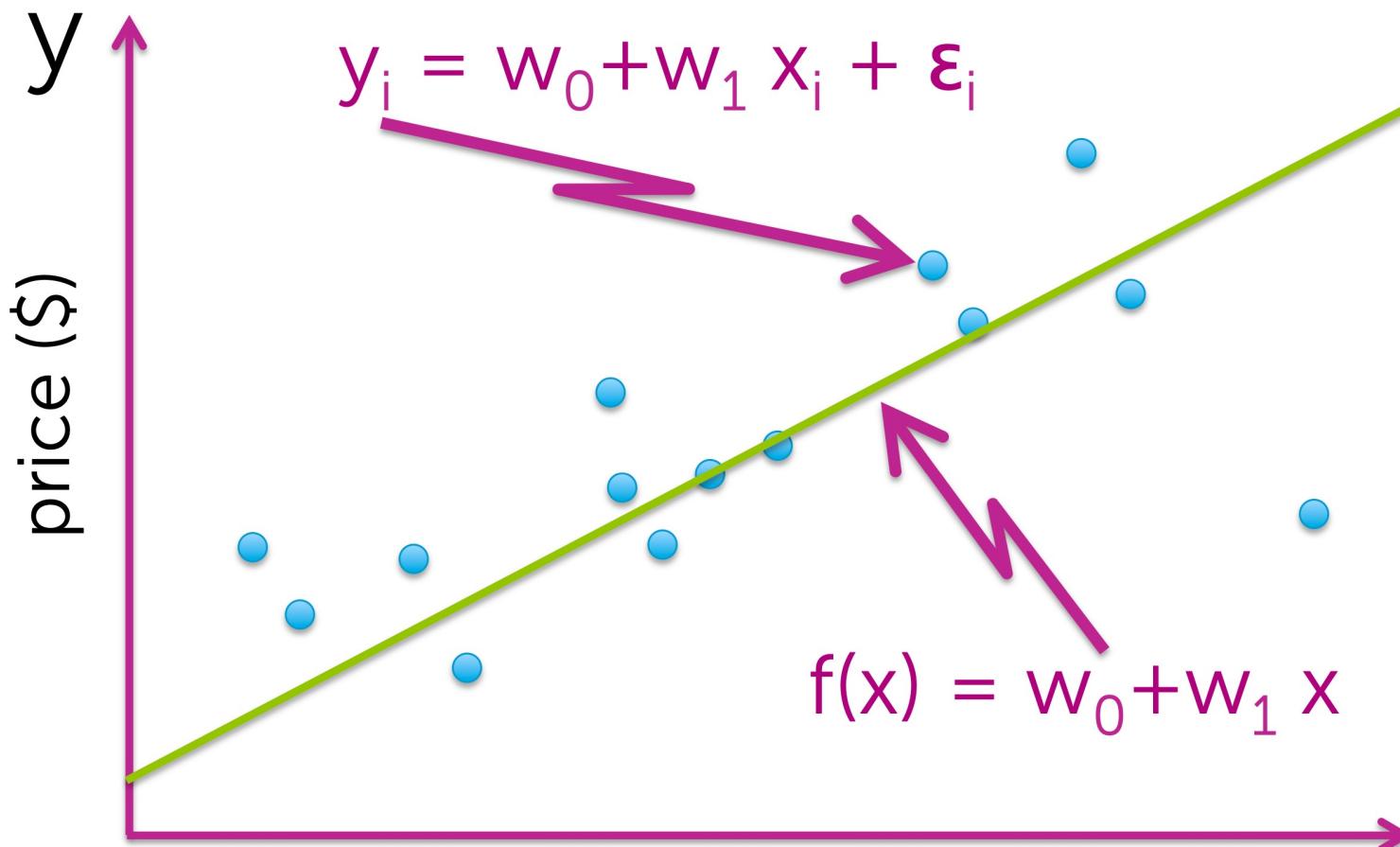
$$\begin{aligned}x_1 &= 4 \\x_2 &= 1 \\y &= 2 \\\hat{y} &= -0.116\end{aligned}$$

Like for another weight

- ...

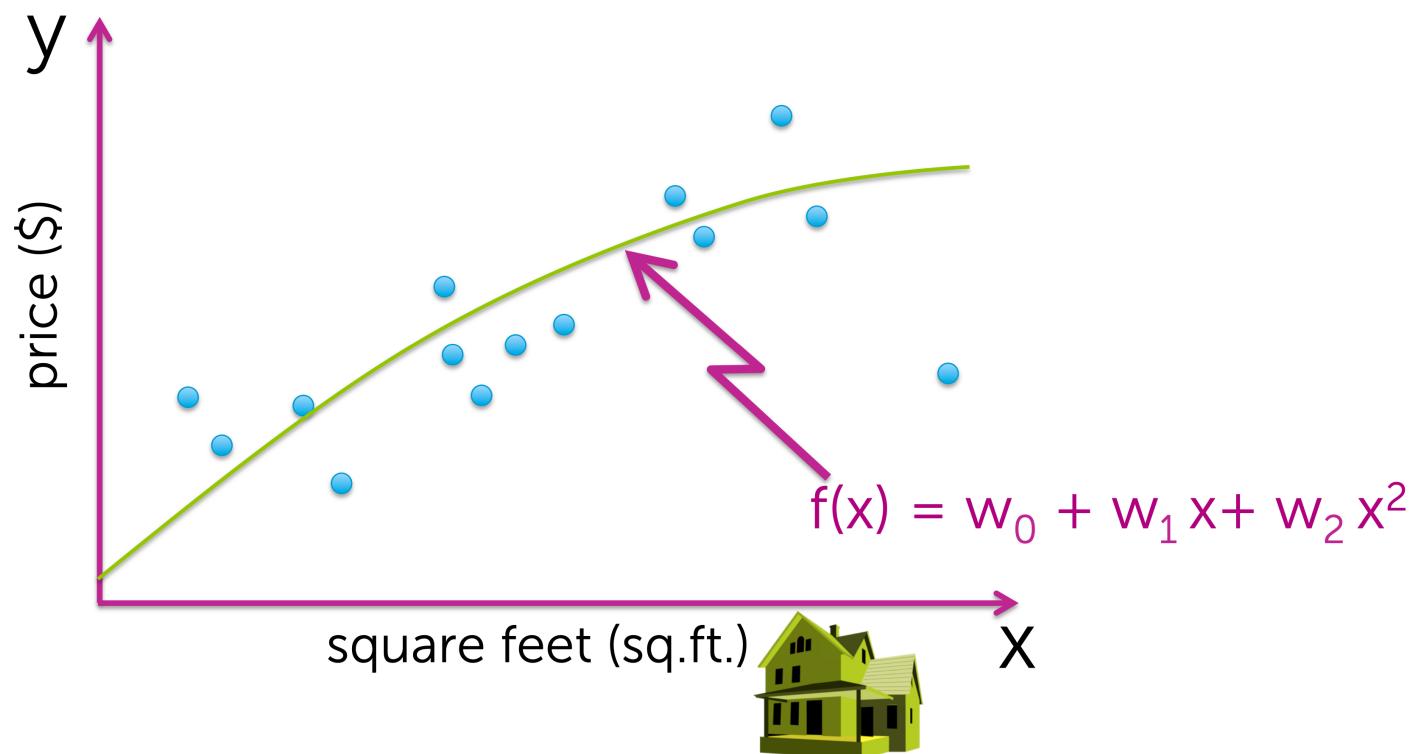
$$\hat{y} = 0.43x_1 - 0.21x_2 + 0.77$$

Simple linear regression model

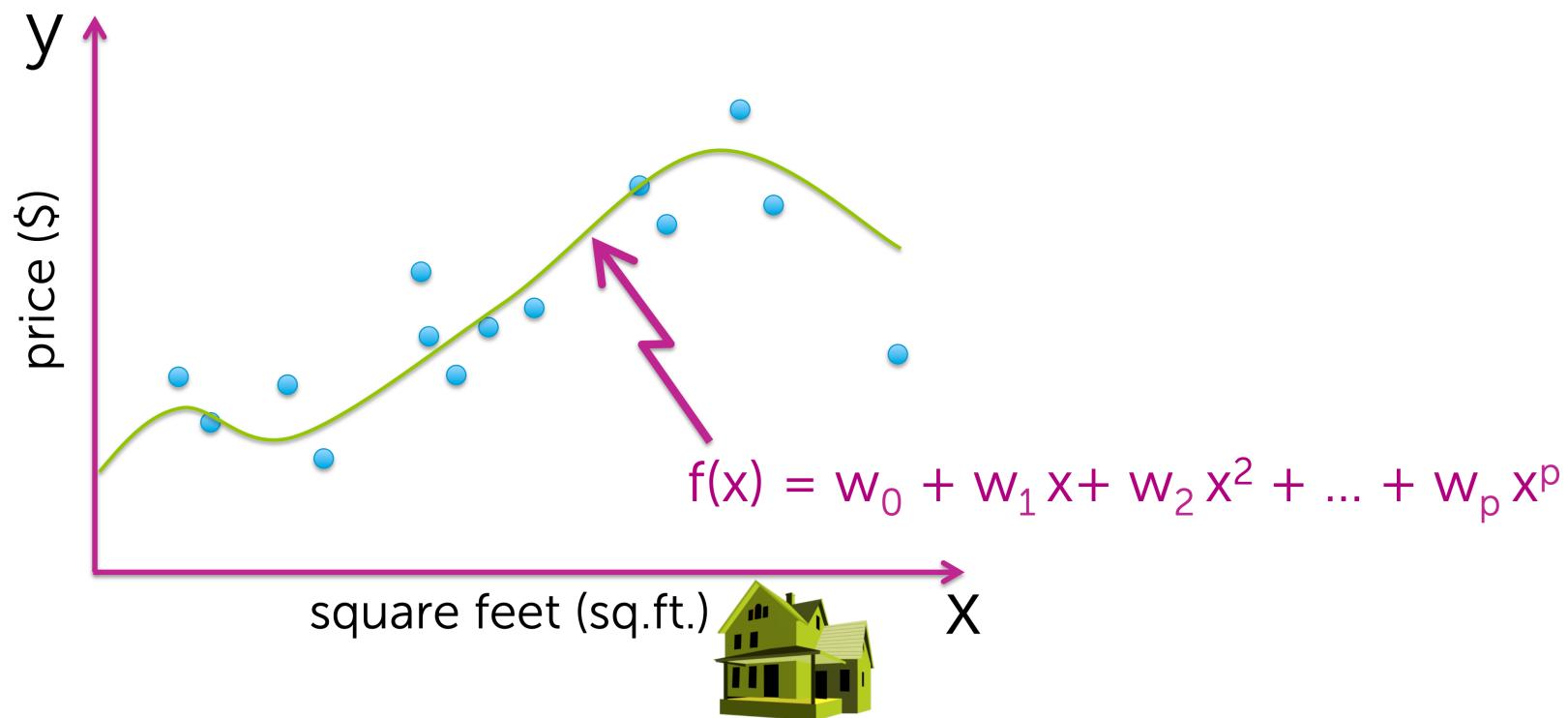


Polynomial regression

What about a quadratic function?



Even higher order polynomial



Polynomial regression

Model:

$$y_i = w_0 + w_1 x_i + w_2 x_i^2 + \dots + w_p x_i^p + \epsilon_i$$

treat as different **features**

feature 1 = 1 (constant) parameter 1 = w₀

feature 3 = x^2 parameter 3 = w_2

■ ■ ■

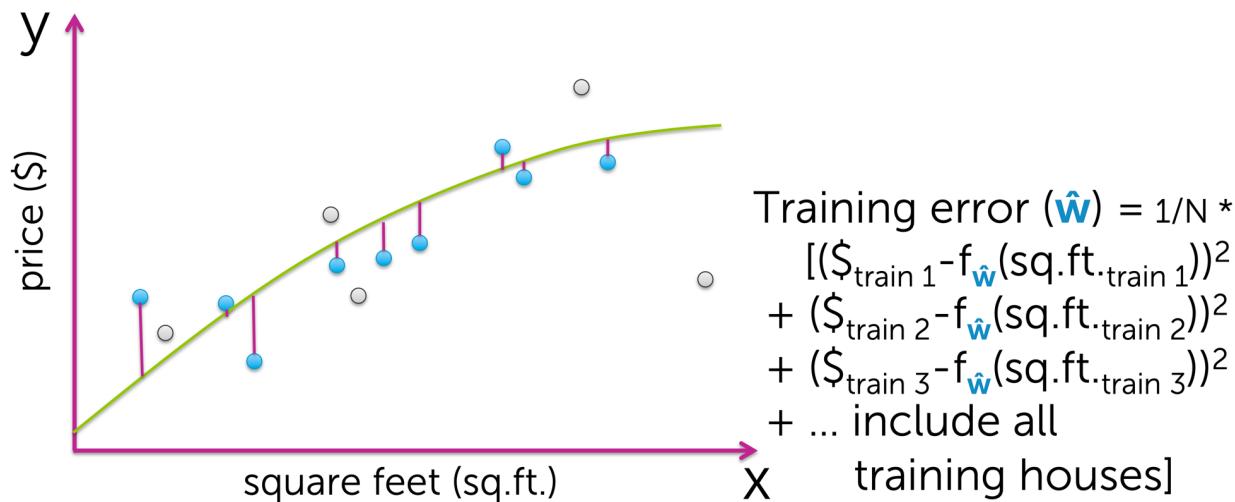
■ ■

feature $p+1 = x^p$

parameter p+1 = w_p

Example:

Use squared error loss $(y - f_{\hat{w}}(x))^2$



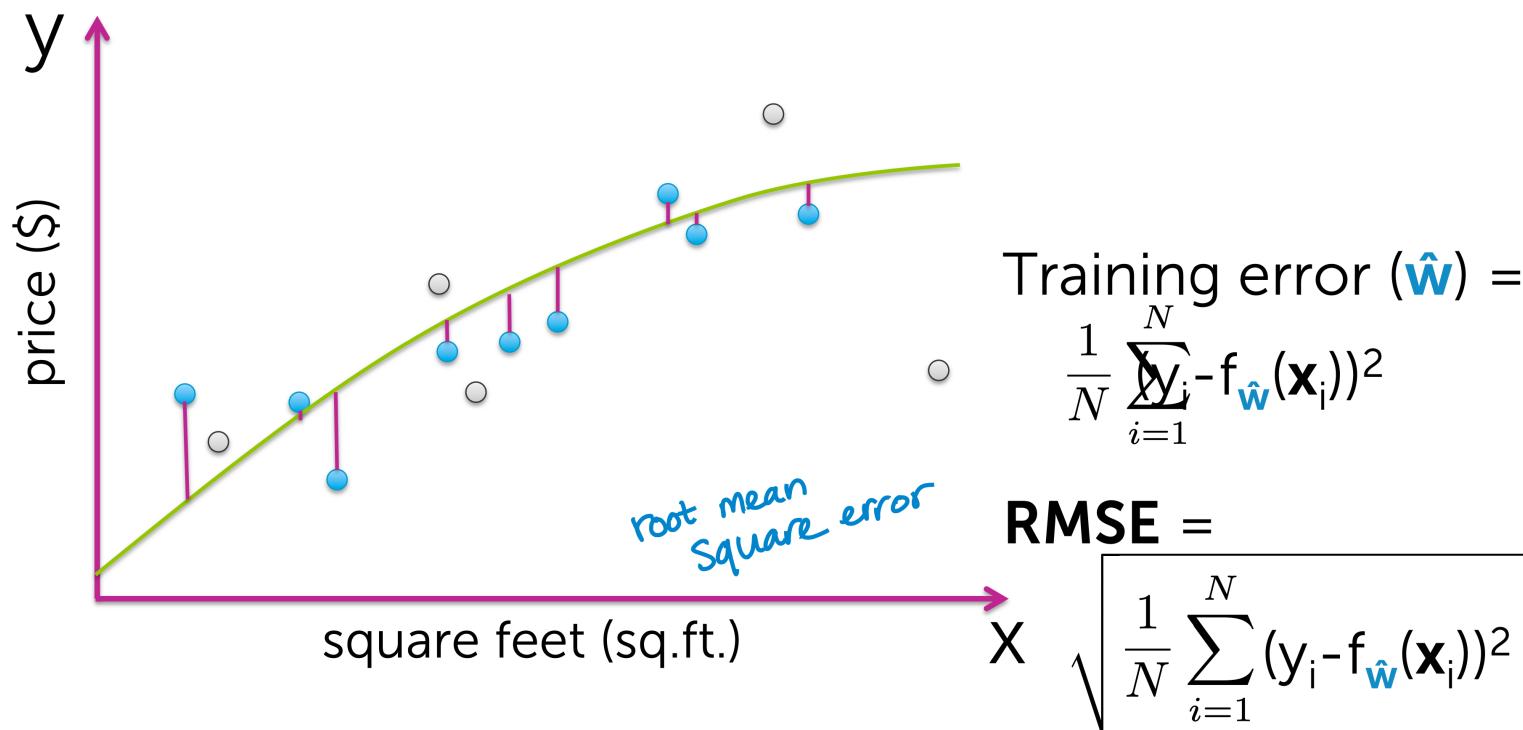
Training error

= avg. loss on houses in training set

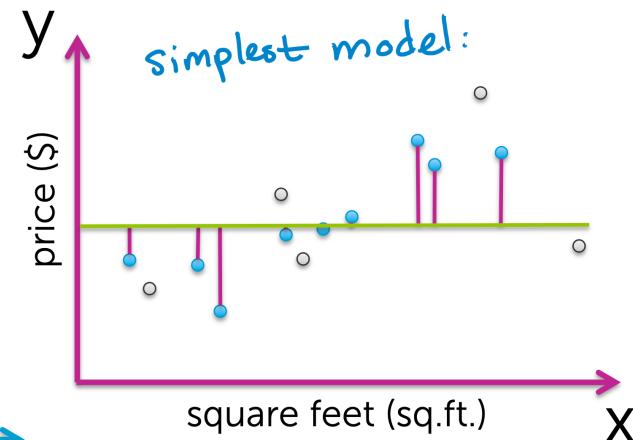
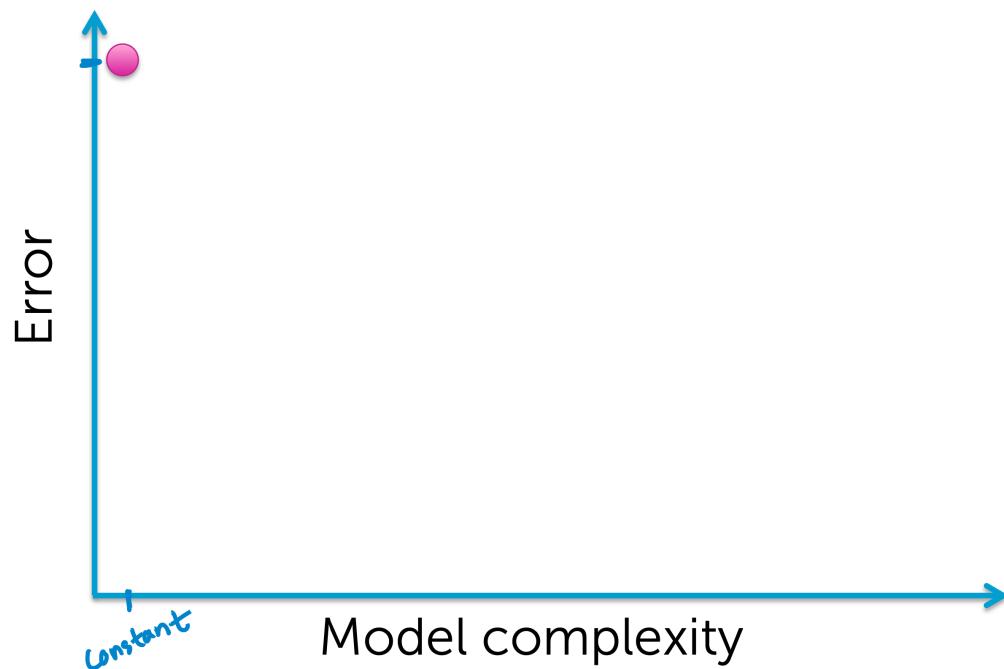
$$= \frac{1}{N} \sum_{i=1}^N L(y_i, f_{\hat{w}}(x_i))$$

fit using training data

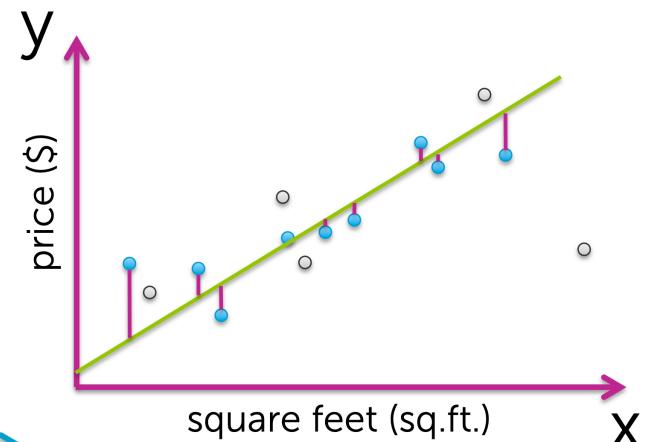
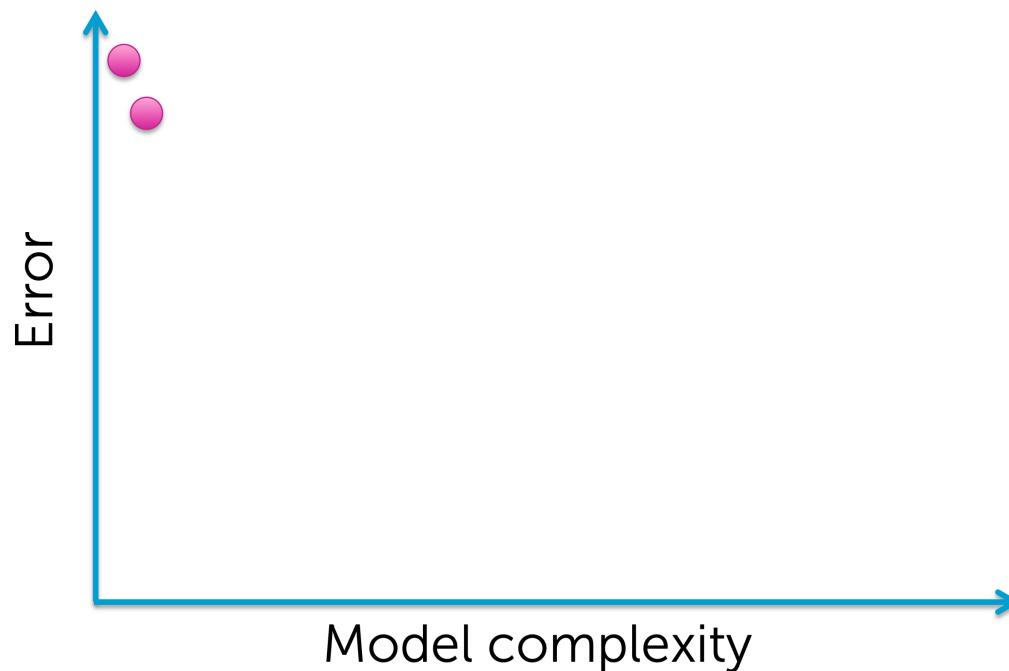
Example:
Use squared error loss $(y - f_{\hat{w}}(x))^2$



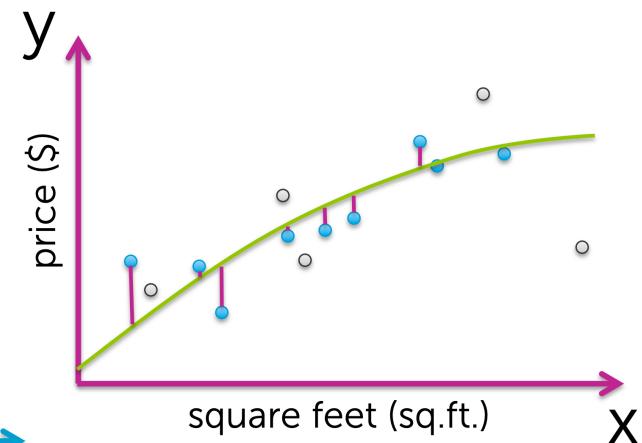
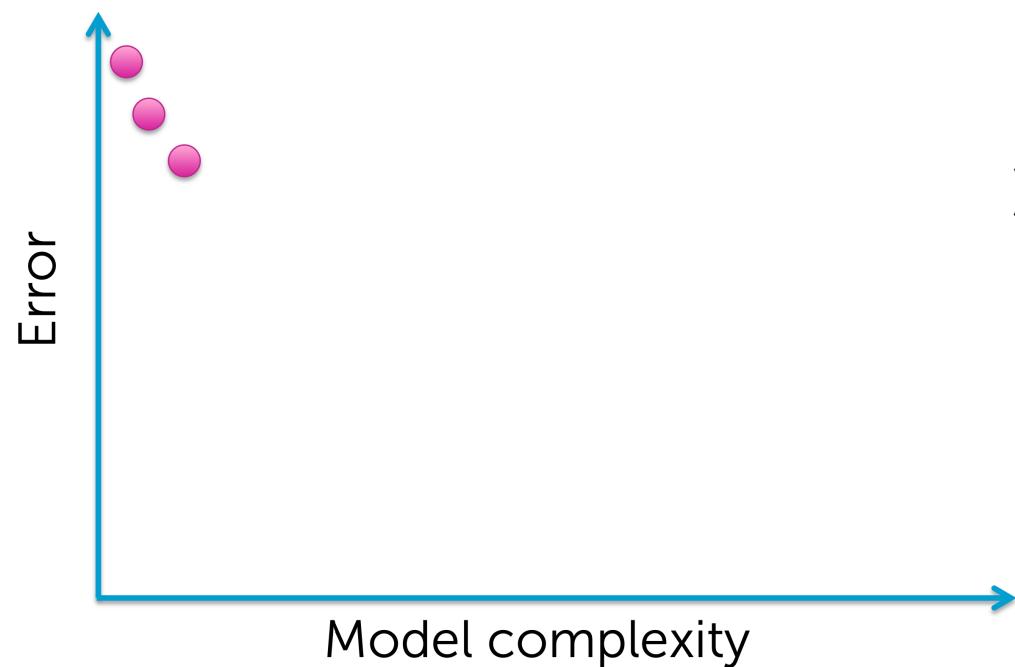
Training error vs. model complexity



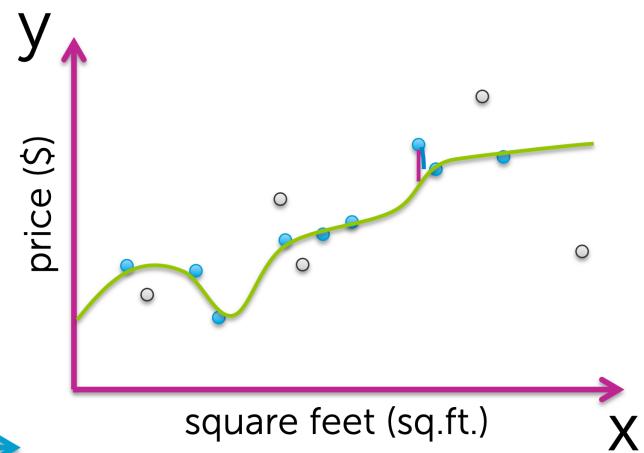
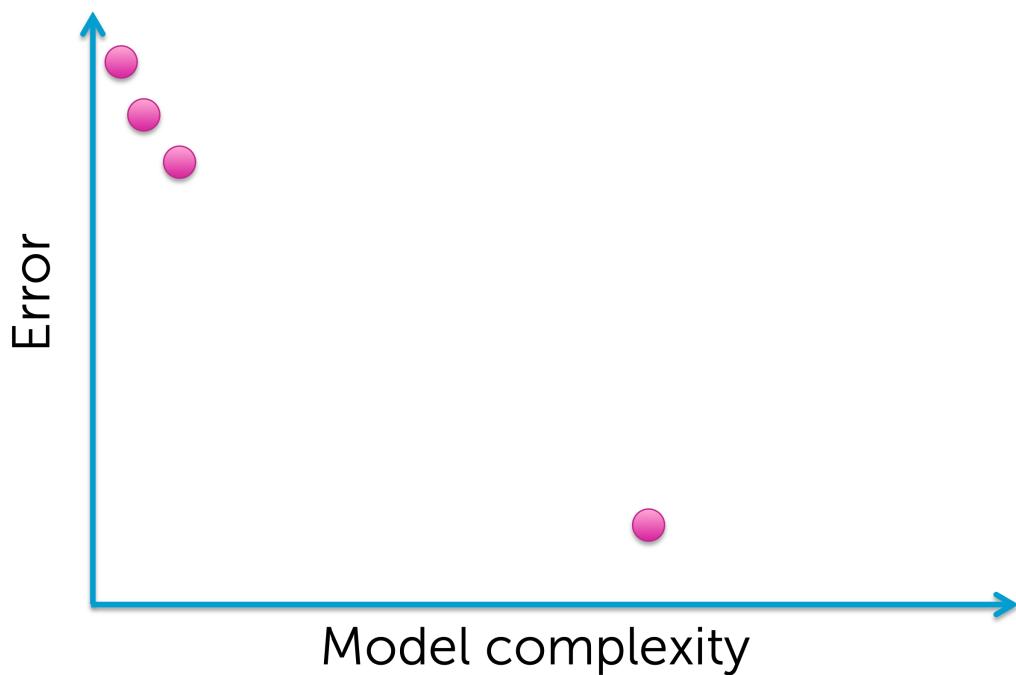
Training error vs. model complexity



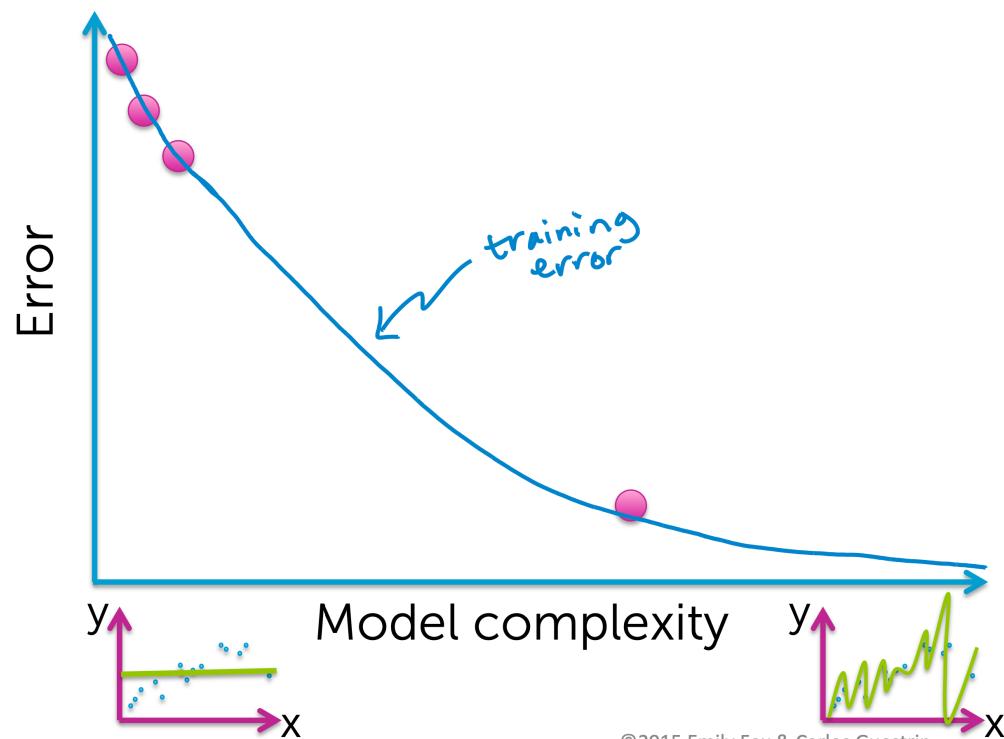
Training error vs. model complexity



Training error vs. model complexity



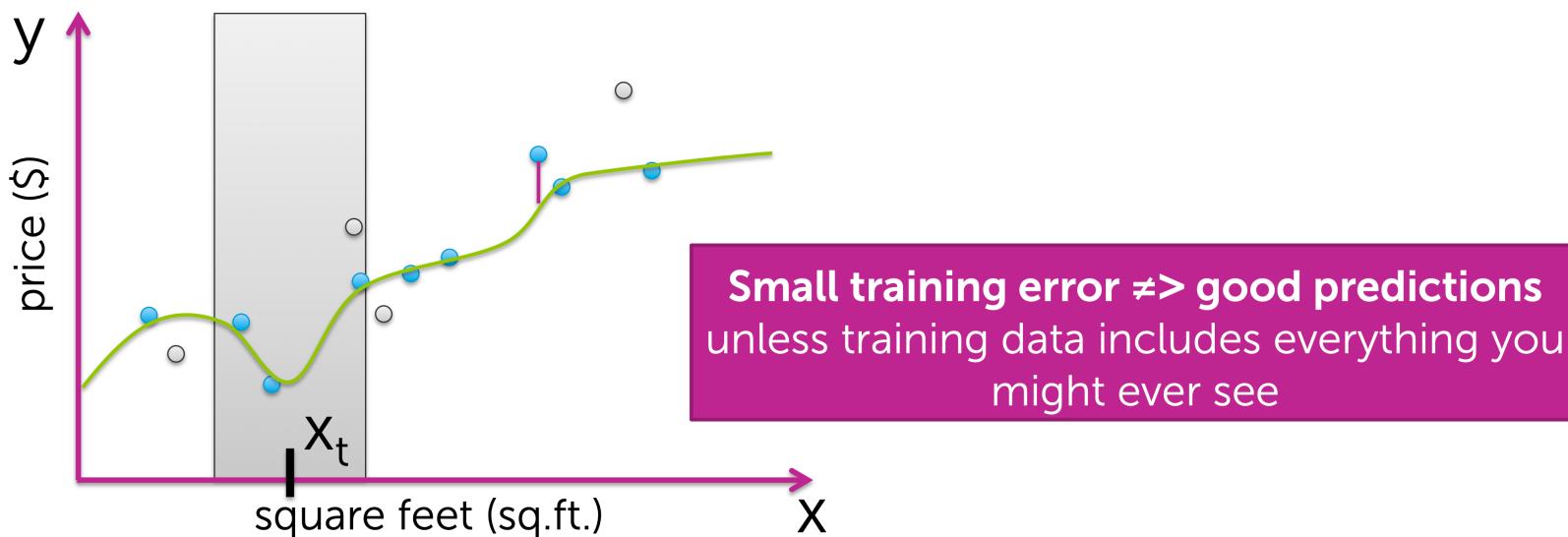
Training error vs. model complexity

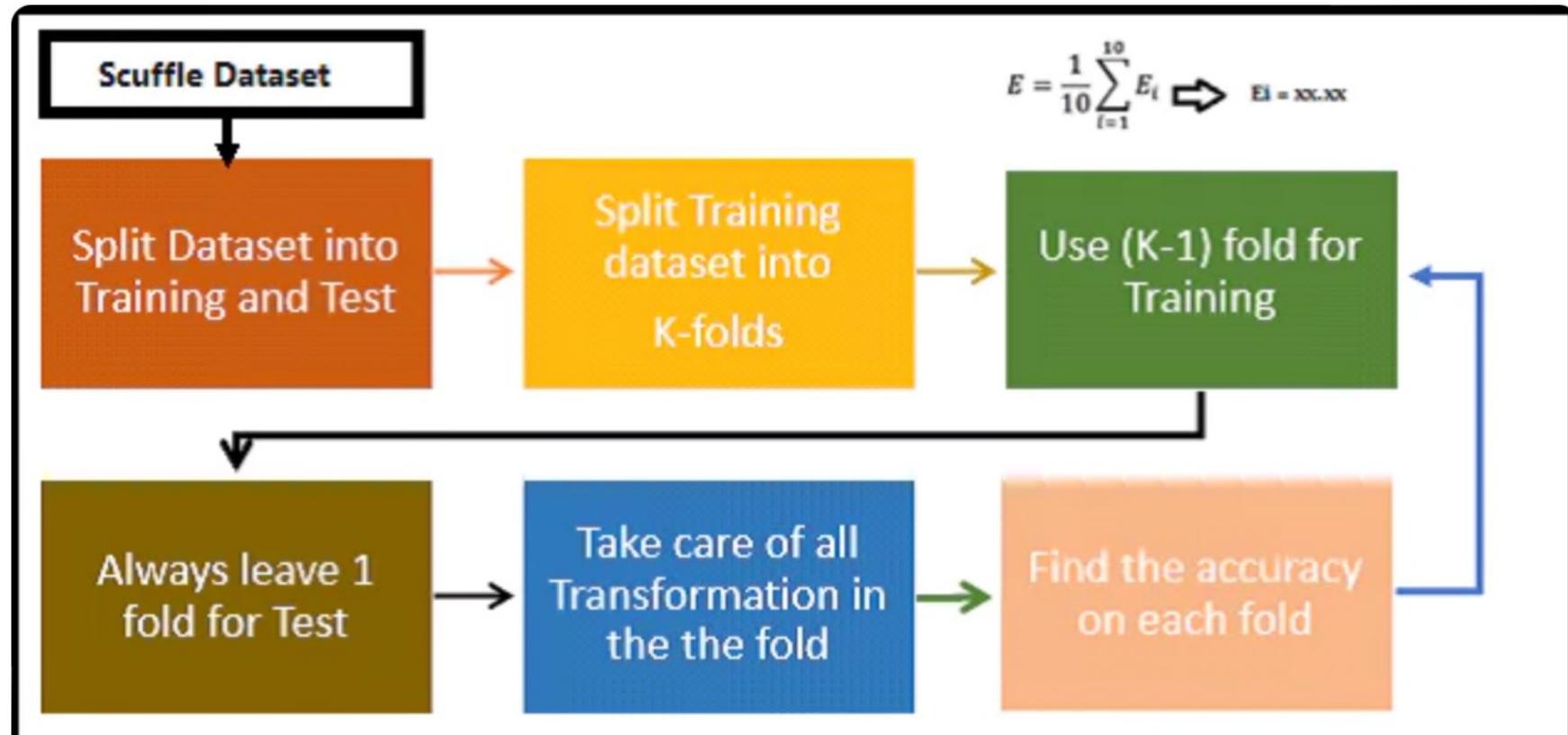


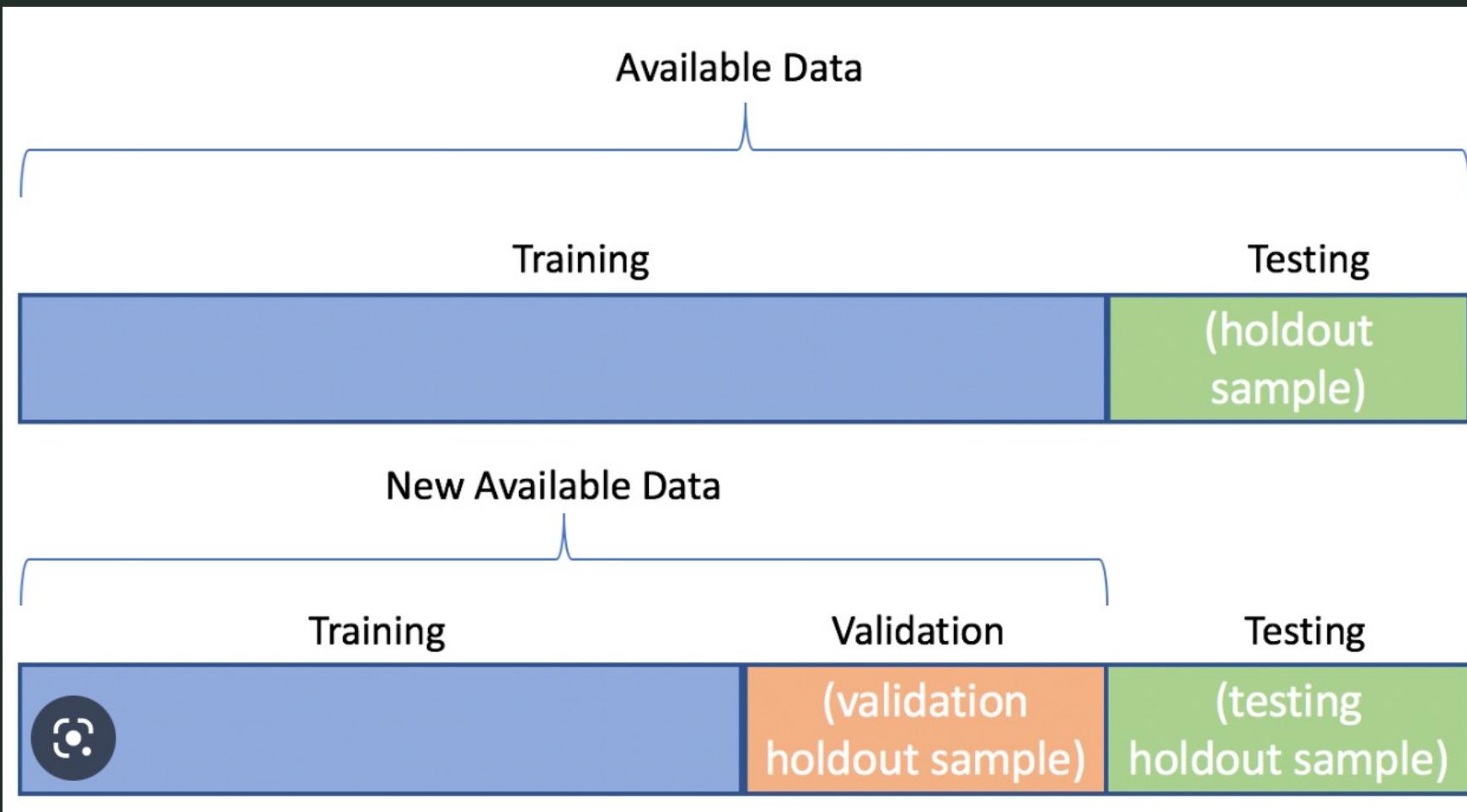
Is training error a good measure of predictive performance?

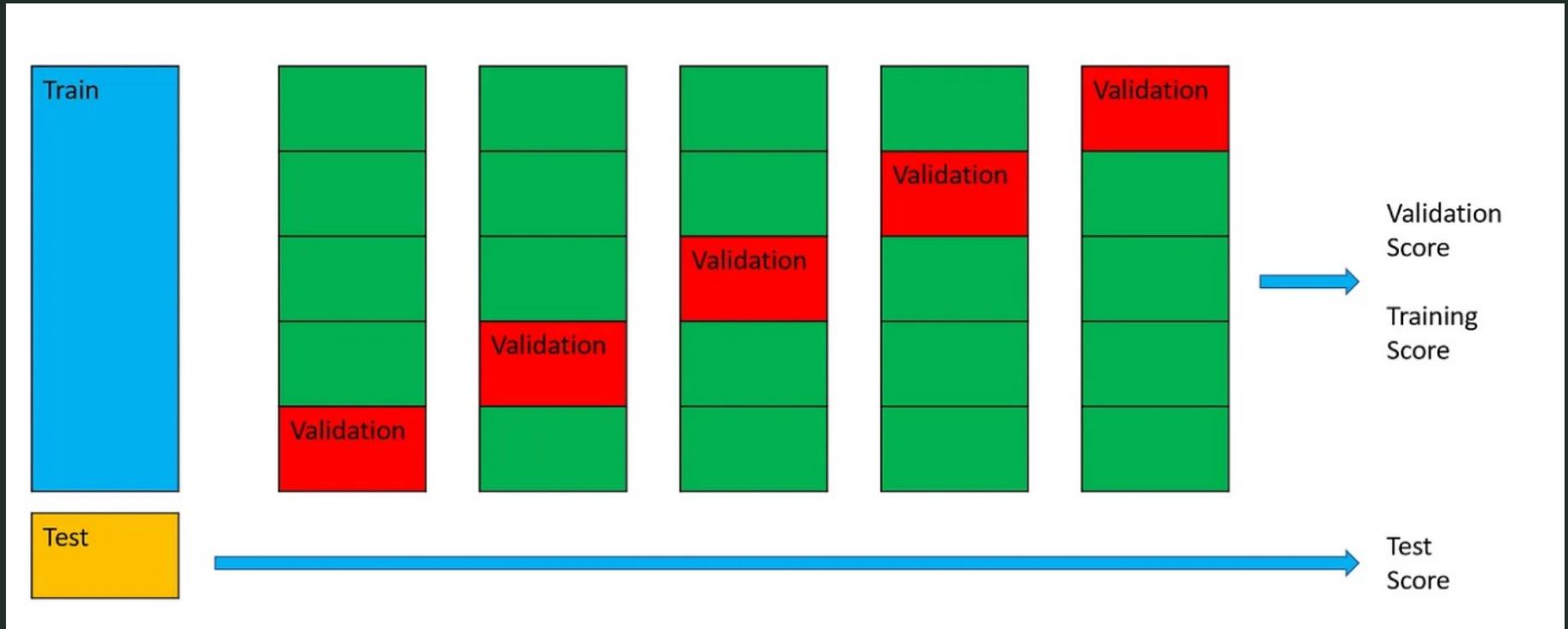
Issue: Training error is overly optimistic

because \hat{w} was fit to training data

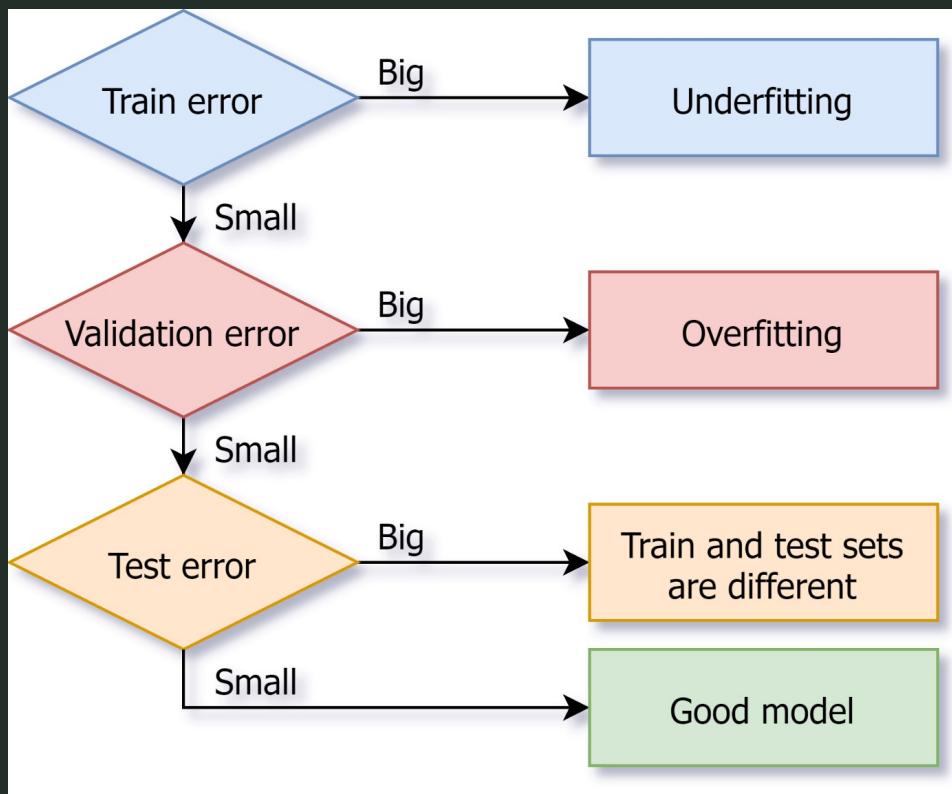


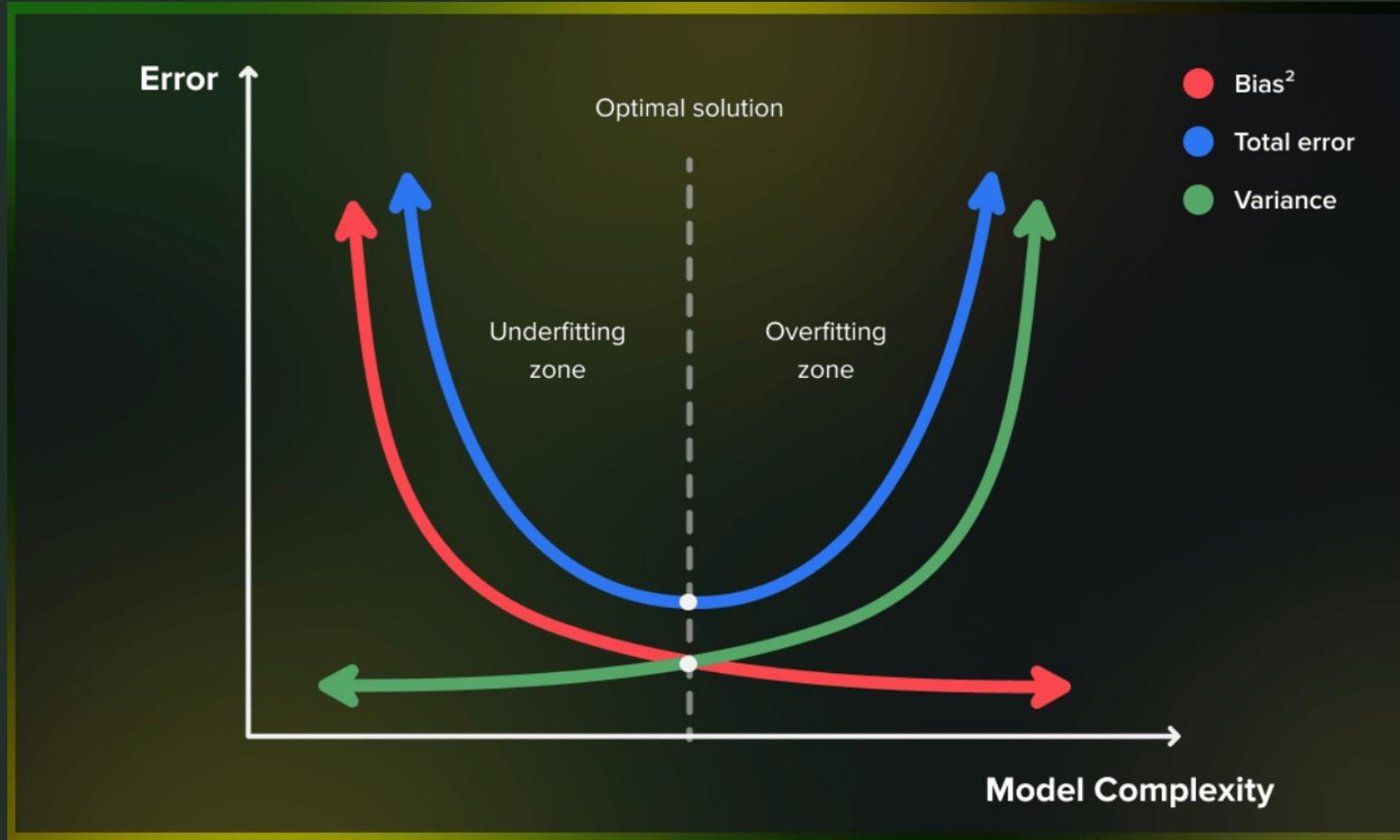






Underfitting Vs Overfitting



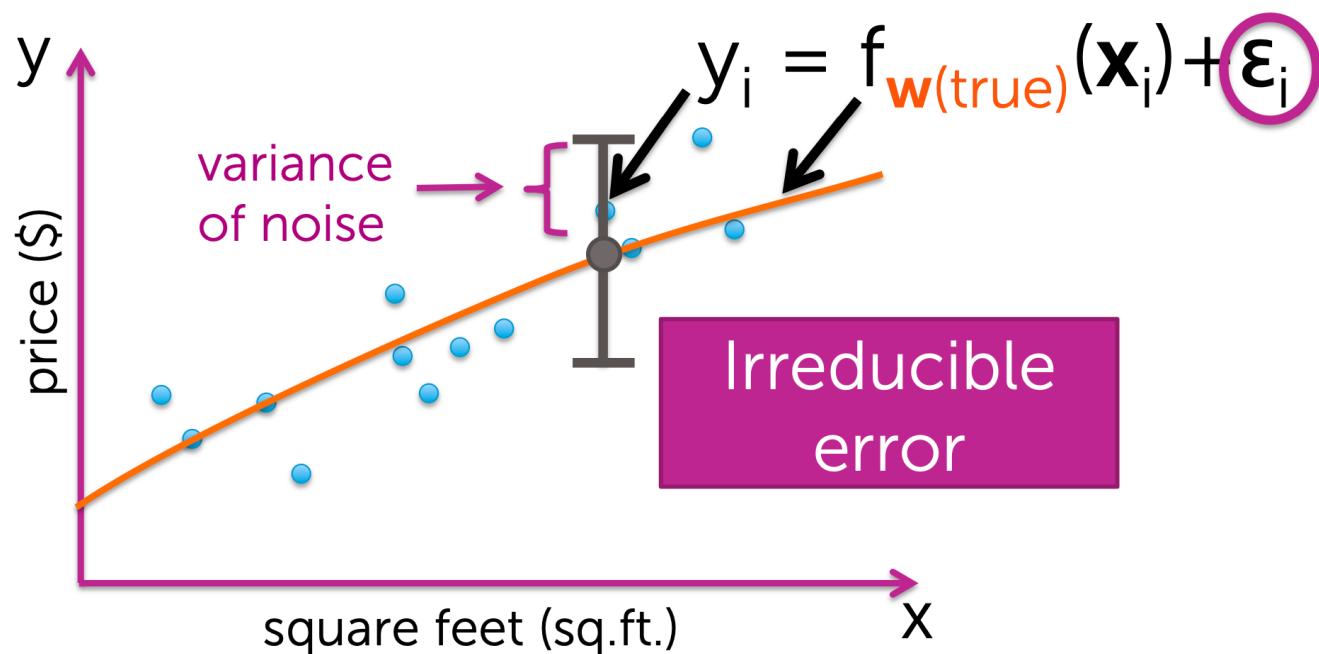


3 sources of error

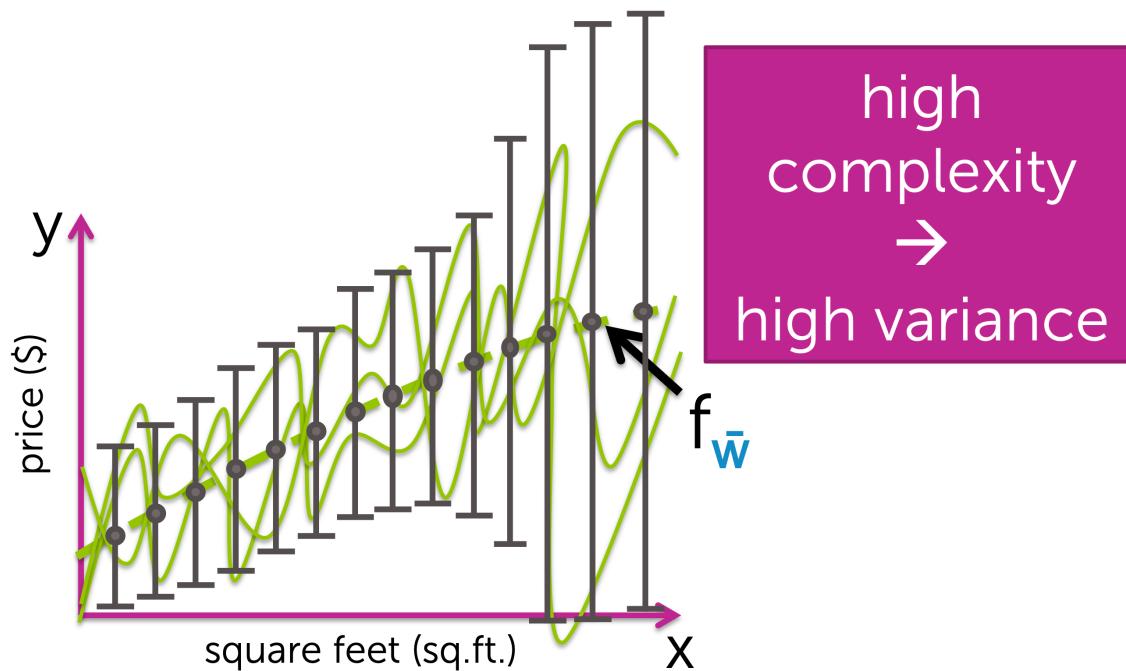
In forming predictions, there are 3 sources of error:

1. Noise
2. Bias
3. Variance

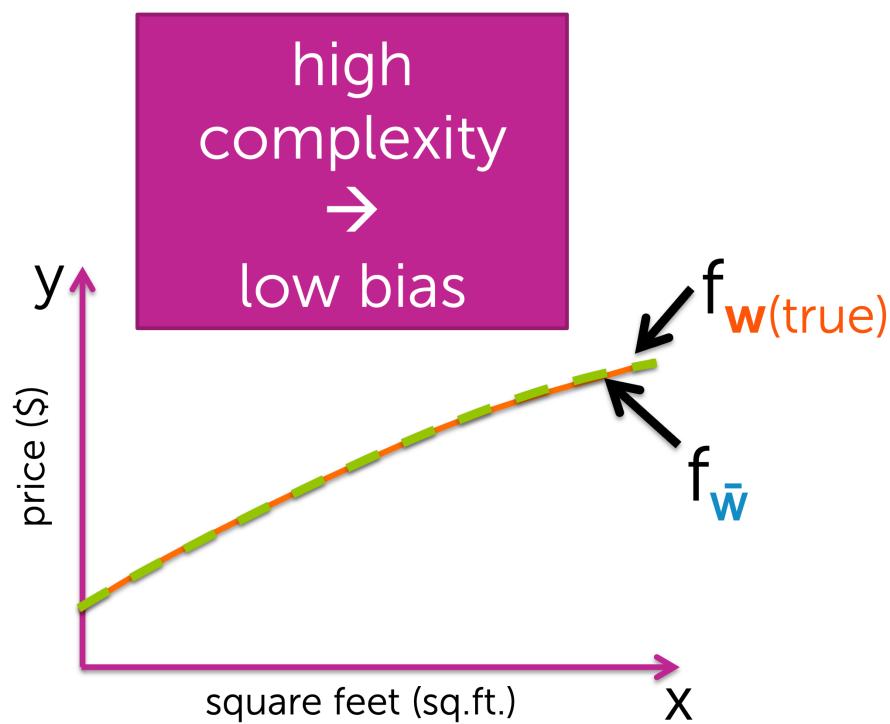
Data inherently noisy



Variance of high-complexity models

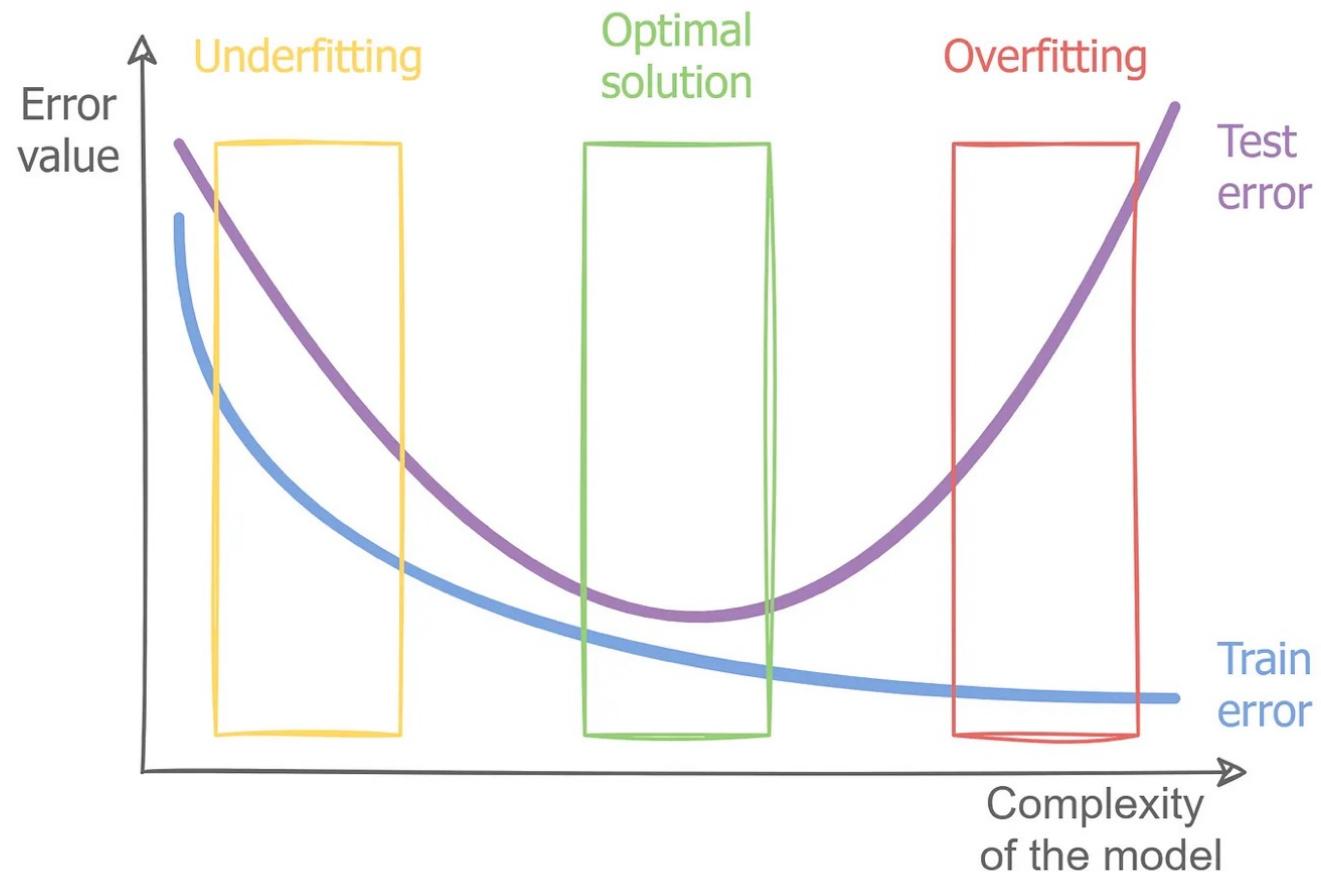


Bias of high-complexity models

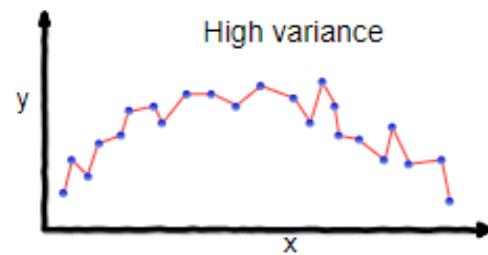




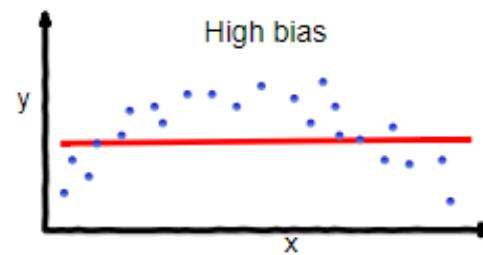
Underfitting Vs Overfitting



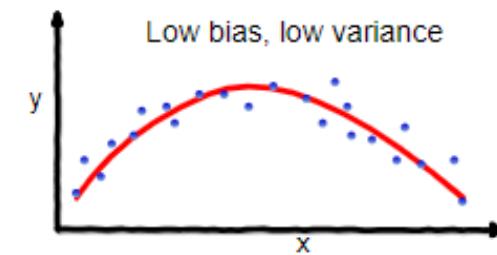
Challenges



overfitting



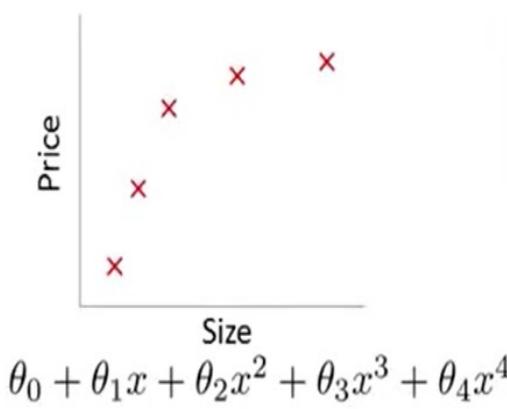
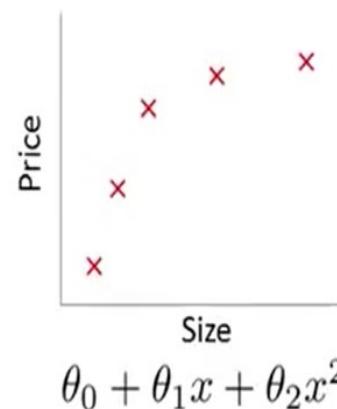
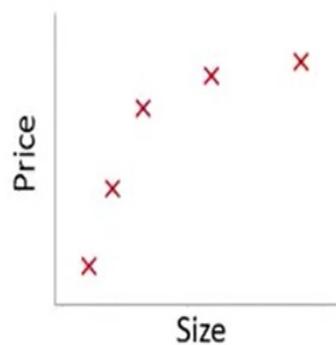
underfitting



Good balance

Underfitting Vs Overfitting in Regression

Example: Linear regression (housing prices)



Solutions

Overfitting and Underfitting are two of the main reasons machine learning models have poor performance.

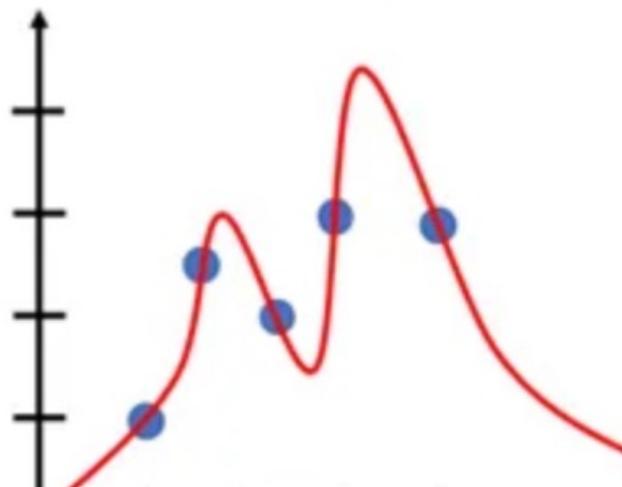
Fixing **overfitting**:

- Simplify the model (fewer parameters)
- Simplify training data (fewer attributes)
- Constrain the model (regularization)
- Use cross-validation
- Use Early stopping
- Build an ensemble
- Gather more data

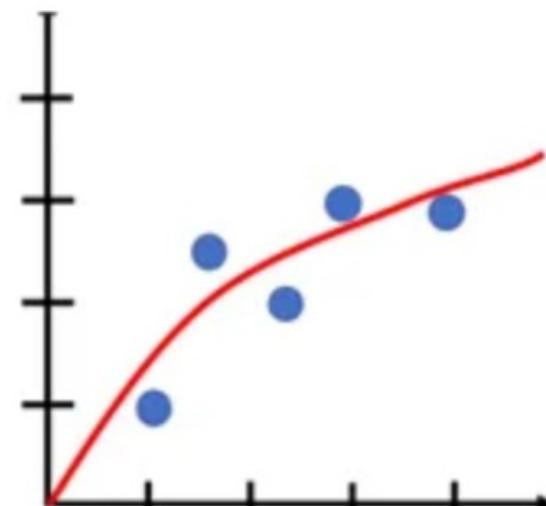
Fixing **underfitting**:

- More complex model (more parameters)
- Increase number of features
- Feature engineer should help
- Un-constrain the model (no regularization)
- Reduce noise on the data
- Train for longer

Impact of Regularization

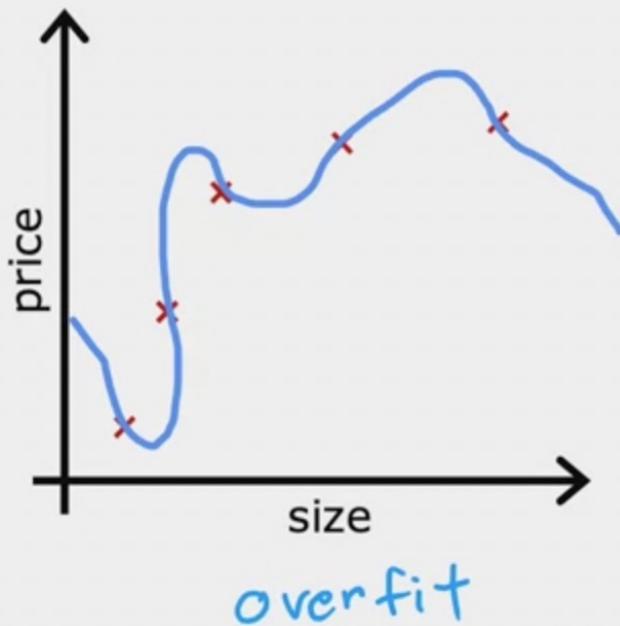


Without Regularization
(Overfit)



With Regularization
(Good fit)

Collect more training examples



Select features to include/exclude

size	bedrooms	floors	age	avg income	...	distance to coffee shop	price
x_1	x_2	x_3	x_4	x_5		x_{100}	y

all features



insufficient data



overfit

selected features

size
bedrooms
age
just right
feature selection

course 2

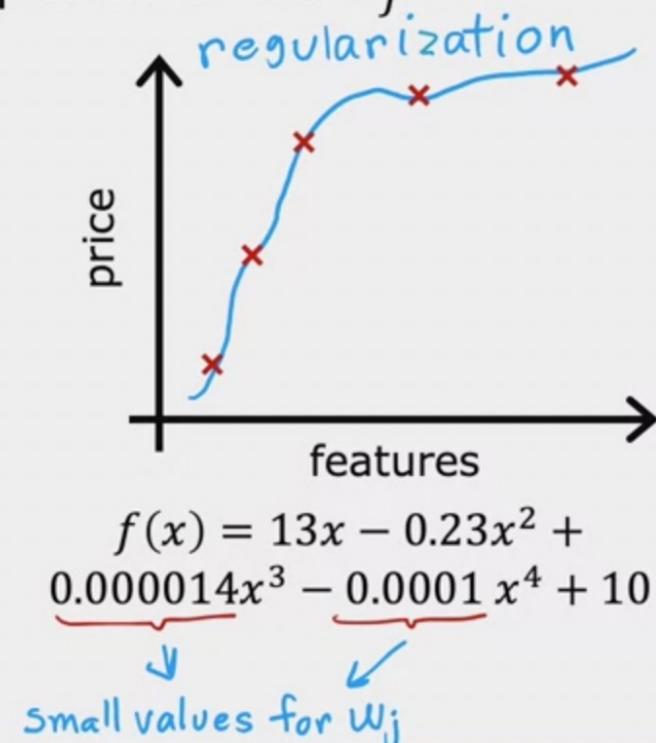
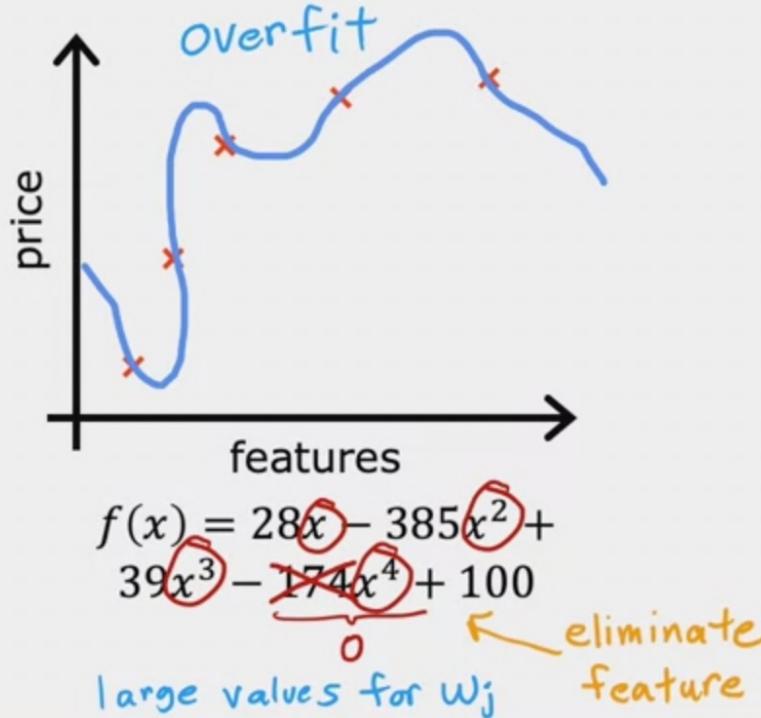
disadvantage



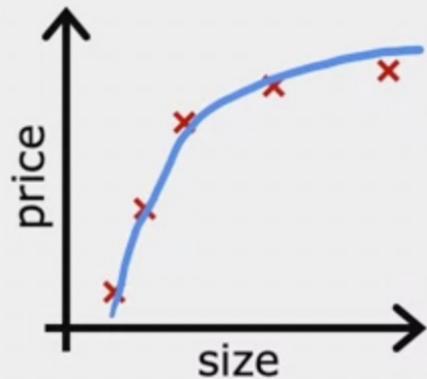
useful features could be lost

Regularization

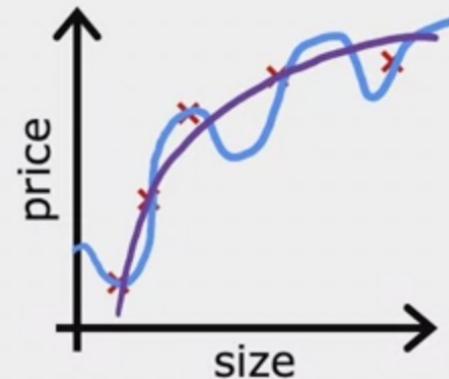
Reduce the size of parameters w_j



Intuition



$$w_1x + w_2x^2 + b$$



$$w_1x + w_2x^2 + w_3x^3 + w_4x^4 + b$$

make w_3, w_4 really small (≈ 0)

$$\min_{\vec{w}, b} \frac{1}{2m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})^2 + 1000 w_3^2 + 1000 w_4^2$$

L1 Regularization(Lasso)

$$\hat{y} = w_1x_1 + w_2x_2 + \dots + w_Nx_N + b$$

$$Loss = Error(y, \hat{y})$$

$$Loss = Error(y, \hat{y}) + \lambda \sum_{i=1}^N |w_i|$$

L2 Regularization(Ridge)

$$\hat{y} = w_1x_1 + w_2x_2 + \dots + w_Nx_N + b$$

$$Loss = Error(y, \hat{y})$$

$$Loss = Error(y, \hat{y}) + \lambda \sum_{i=1}^N w_i^2$$

Elastic Net Regularization

$$Loss = Error(y, \hat{y}) + \alpha \cdot \lambda \sum_{i=1}^N |w_i| + (1-\alpha) \lambda \sum_{i=1}^N w_i^2$$