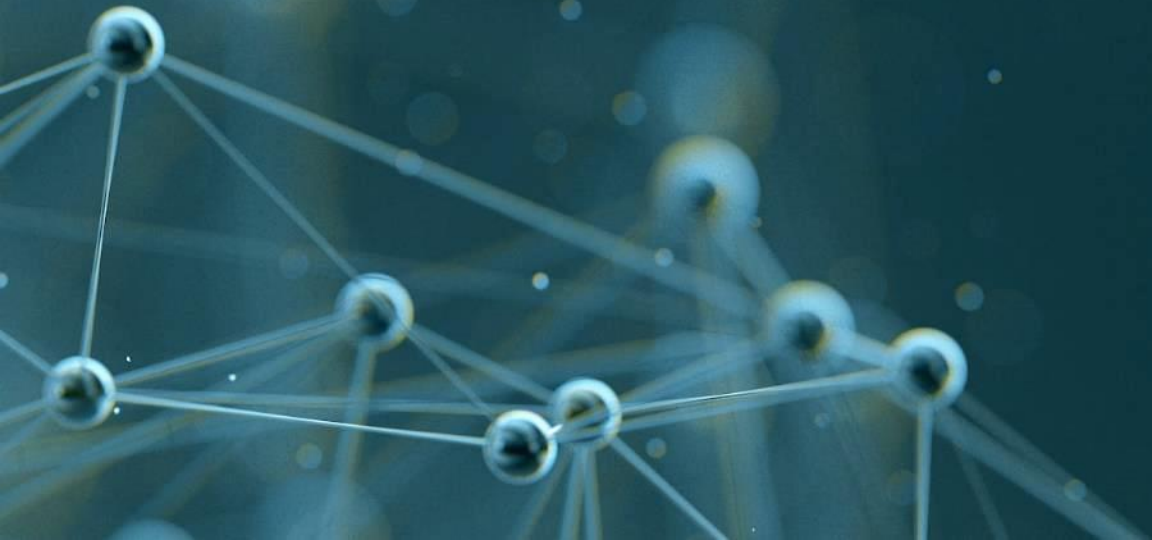
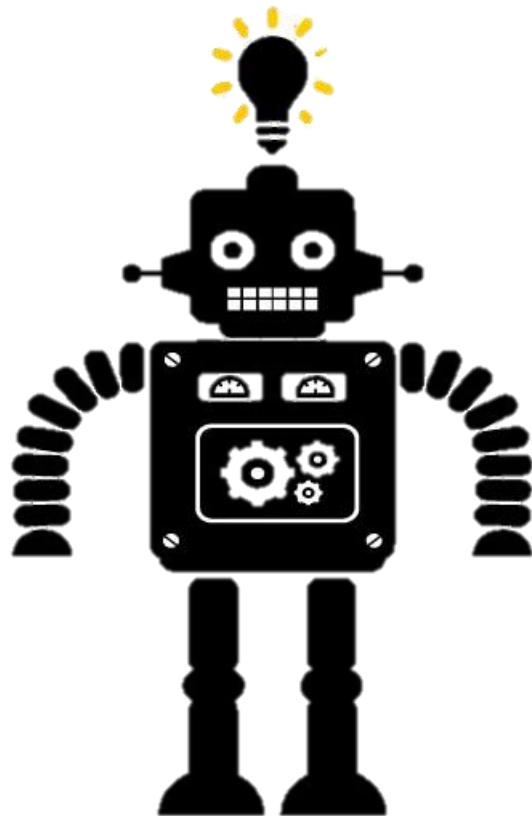


REVIEW OF MACHINE LEARNING



WHAT IS MACHINE LEARNING?

Machine learning allows computers to learn and infer from data.



TYPES OF MACHINE LEARNING

SUPERVISED

Data points have known outcome

UNSUPERVISED

Data points have unknown outcome

TYPES OF SUPERVISED LEARNING

REGRESSION

Outcome is continuous (numerical)

CLASSIFICATION

Outcome is a category

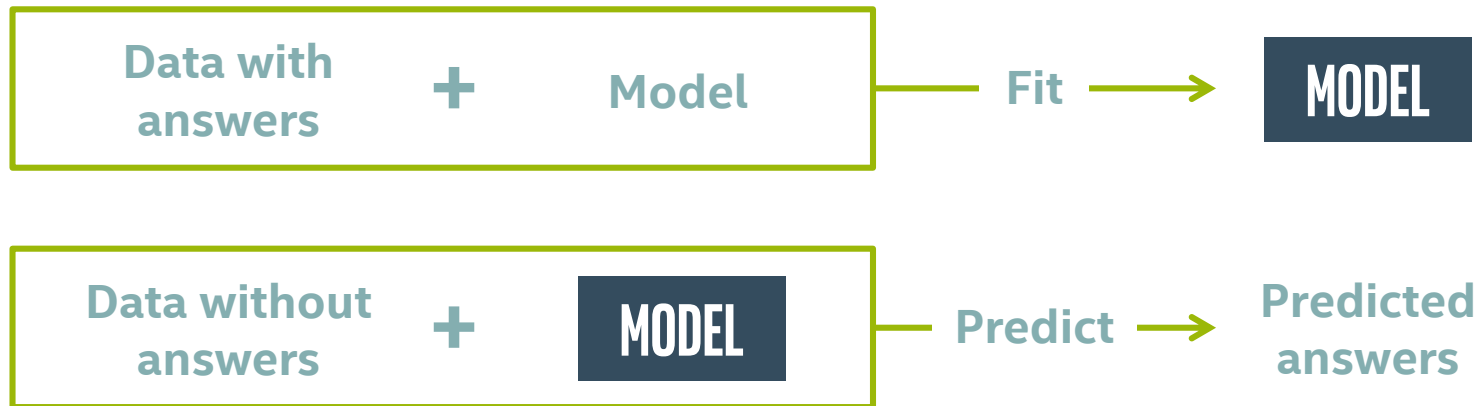
MACHINE LEARNING VOCABULARY

- **Target:** predicted category or value of the data
(column to predict)
- **Features:** properties of the data used for prediction
(non-target columns)
- **Example:** a single data point within the data
(one row)
- **Label:** the target value for a single data point

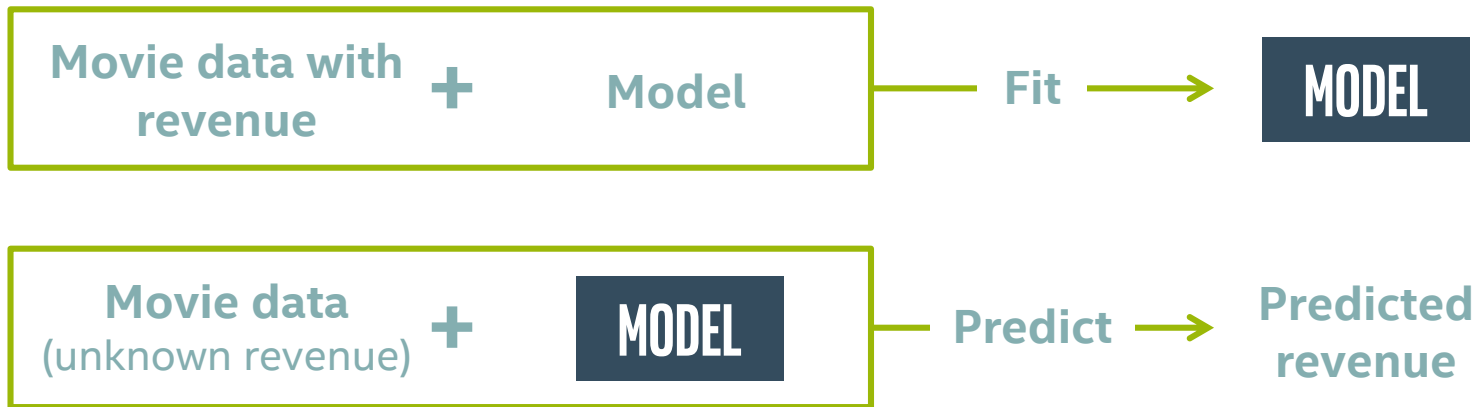
MACHINE LEARNING VOCABULARY (SYNONYMS)

- **Target:** Response, Output, Dependent Variable, Labels
- **Features:** Predictors, Input, Independent Variables, Attributes
- **Example:** Observation, Record, Instance, Datapoint, Row
- **Label:** Answer, y-value, Category

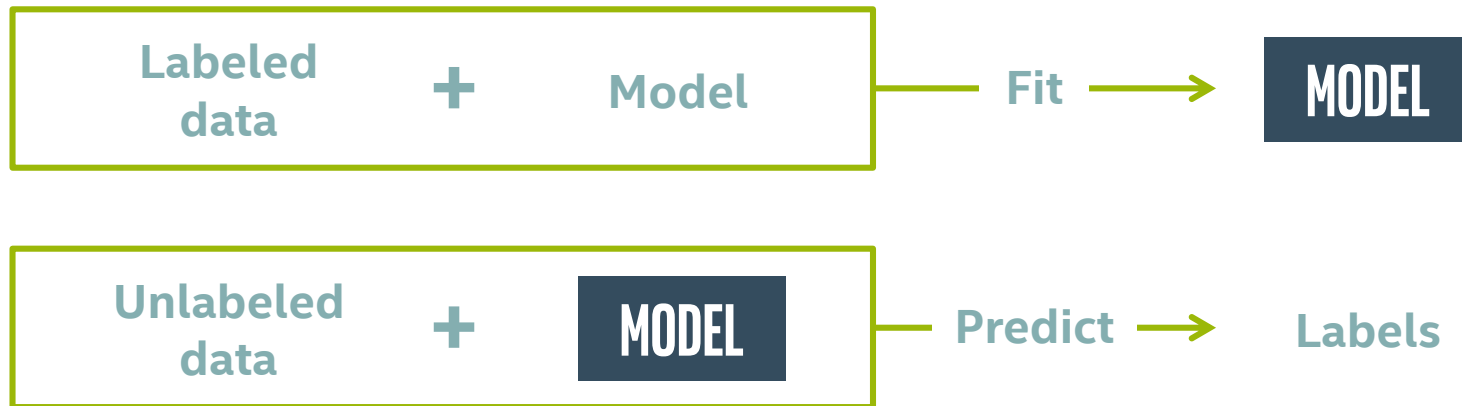
SUPERVISED LEARNING OVERVIEW



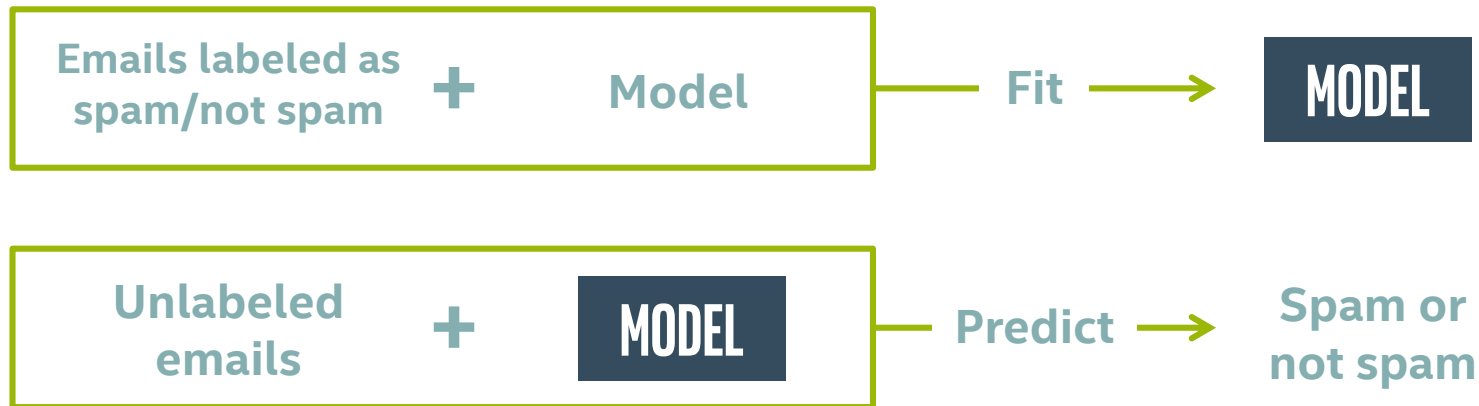
REGRESSION: NUMERICAL ANSWERS



CLASSIFICATION: CATEGORICAL ANSWERS



CLASSIFICATION: CATEGORICAL ANSWERS



THREE TYPES OF CLASSIFICATION PREDICTIONS

- **Hard Prediction:** Predict a single category for each instance.
- **Ranking Prediction:** Rank the instances from most likely to least likely. (binary classification)
- **Probability Prediction:** Assign a probability distribution across the classes to each instance.

METRICS FOR CLASSIFICATION

- **Hard Prediction:** Accuracy, Precision, Recall (Sensitivity), Specificity, F1 Score
- **Ranking Prediction:** AUC (ROC), Precision-Recall Curves
- **Probability Prediction:** Log-loss (aka Cross-Entropy), Brier Score

METRICS FOR REGRESSION

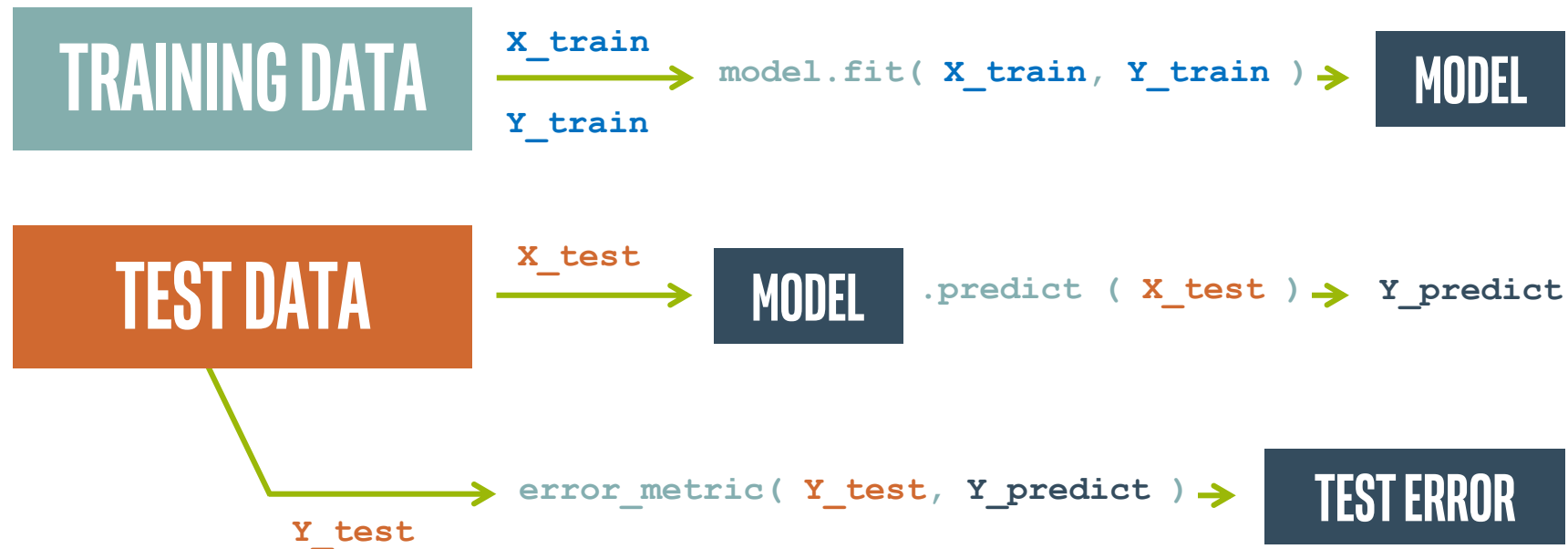
Root Mean Square Error (RMSE)

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$$

Mean Absolute Deviation

$$MAD = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

FITTING TRAINING AND TEST DATA



USING TRAINING AND TEST DATA

TRAINING DATA

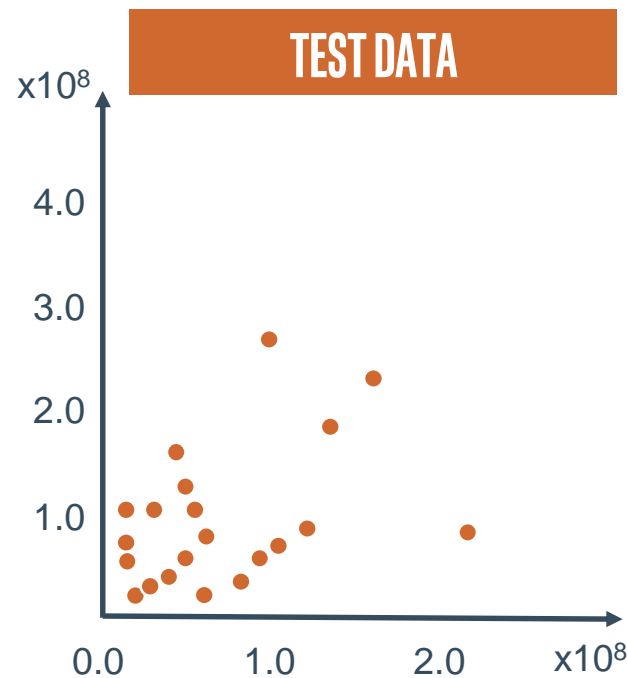
fit the model

TEST DATA

measure performance

- predict label with model
- compare with actual value
- measure error

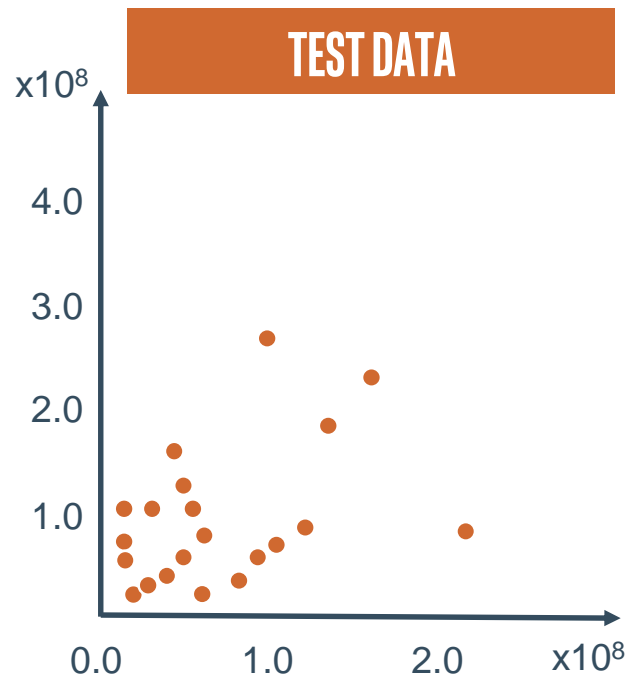
USING TRAINING AND TEST DATA



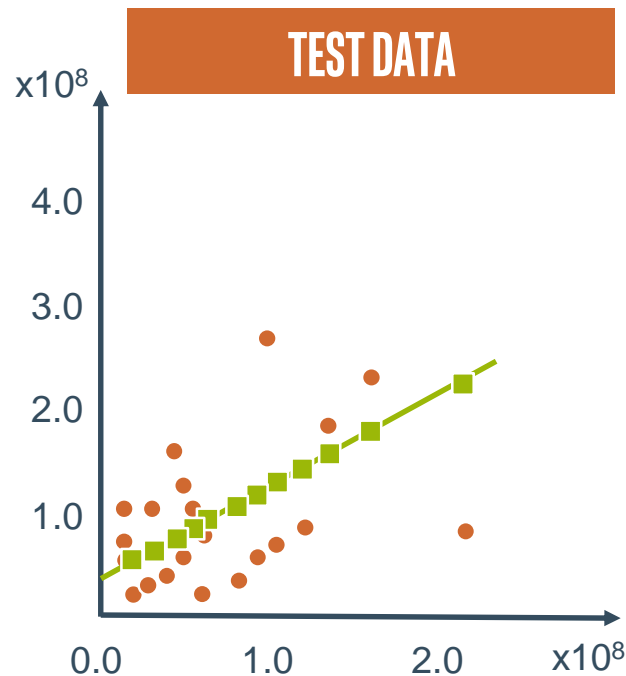
USING TRAINING AND TEST DATA



Fit the model

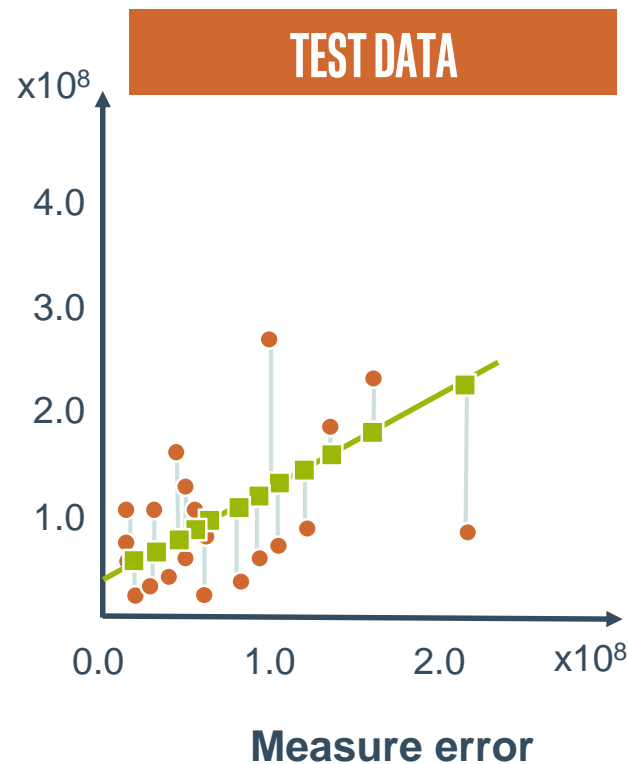


USING TRAINING AND TEST DATA



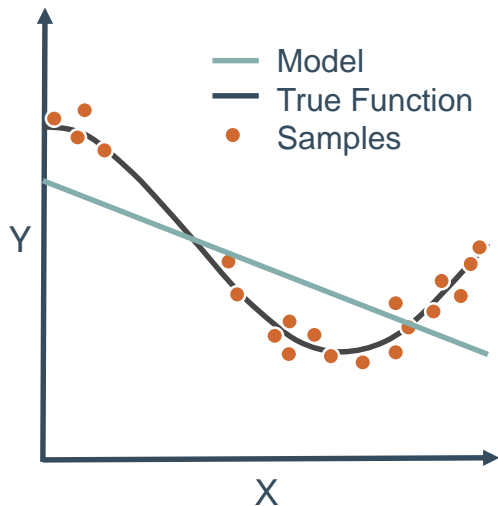
Make predictions

USING TRAINING AND TEST DATA



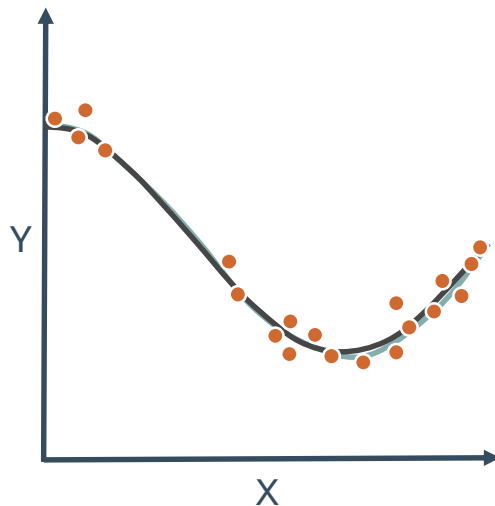
HOW WELL DOES THE MODEL GENERALIZE?

Polynomial Degree = 1



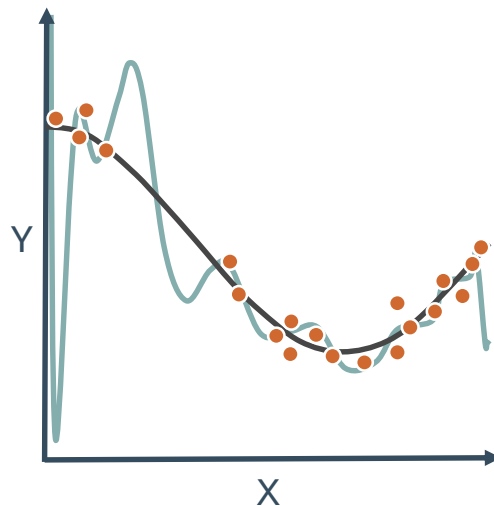
Poor at Training Set
Poor at Predicting

Polynomial Degree = 4



Just Right

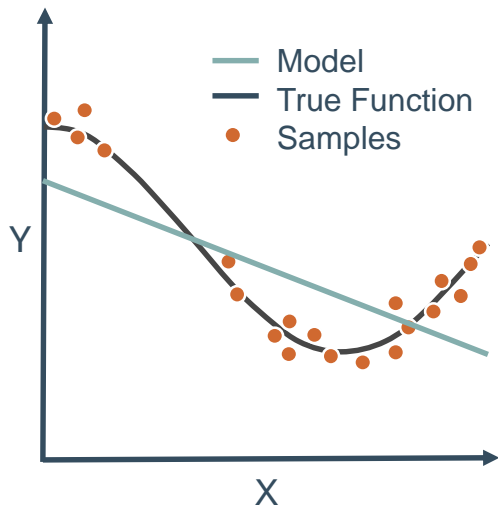
Polynomial Degree = 15



Good at Training Set
Poor at Predicting

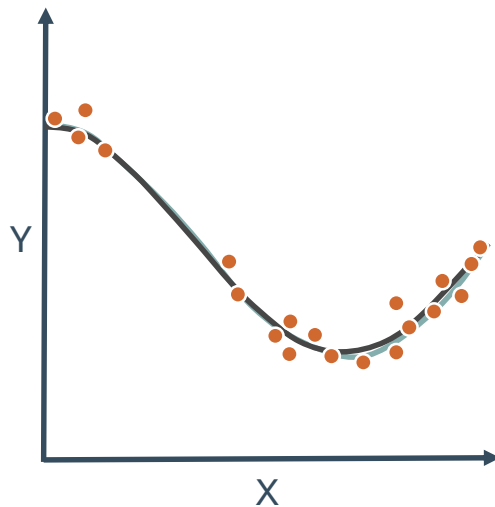
UNDERFITTING VS OVERFITTING

Polynomial Degree = 1



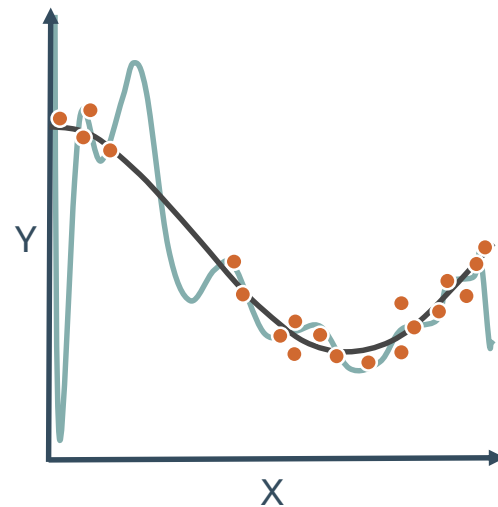
Underfitting

Polynomial Degree = 4



Just Right

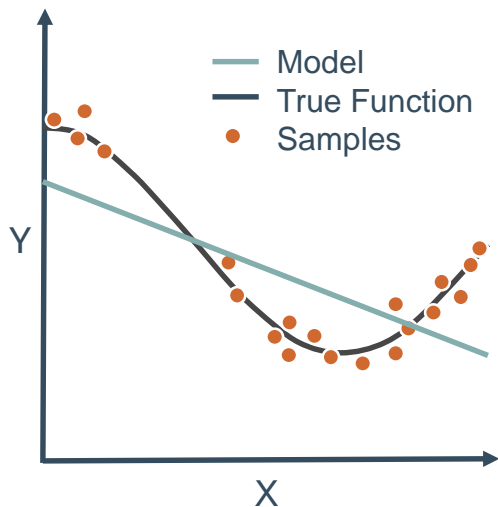
Polynomial Degree = 15



Overfitting

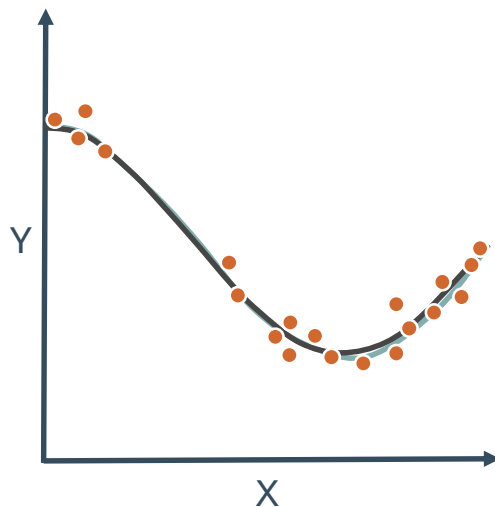
BIAS—VARIANCE TRADEOFF

Polynomial Degree = 1



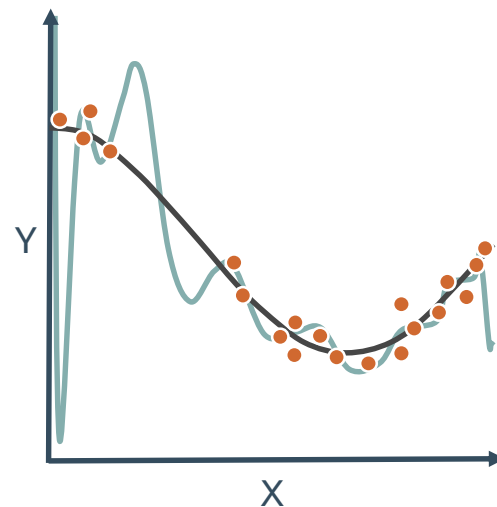
**High Bias
Low Variance**

Polynomial Degree = 4



Just Right

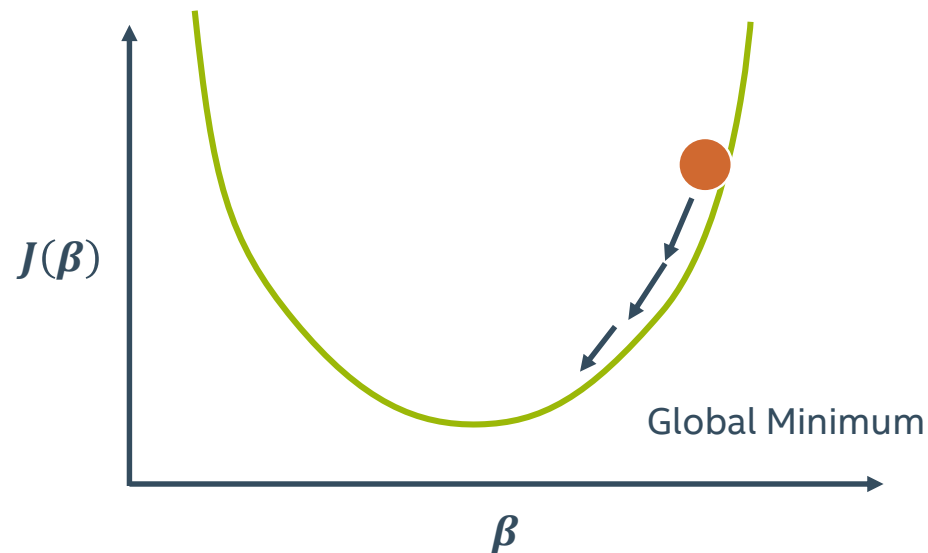
Polynomial Degree = 15



**Low Bias
High Variance**

GRADIENT DESCENT

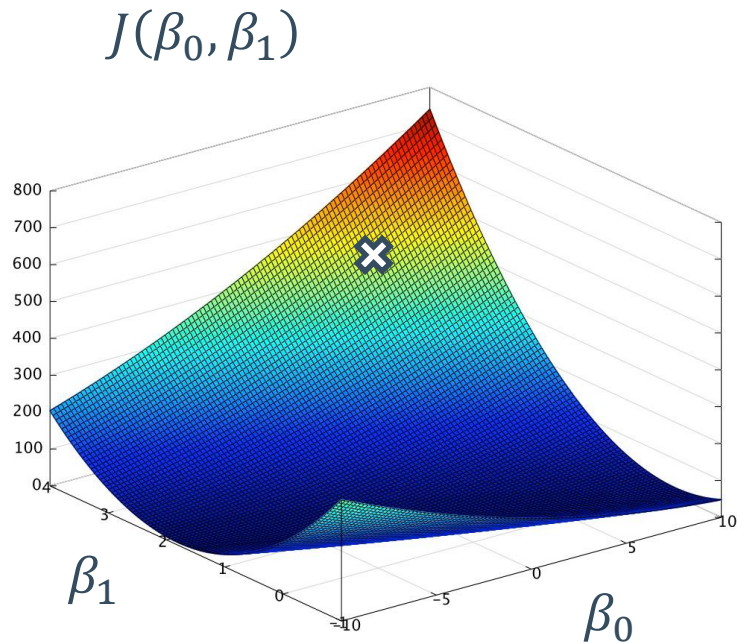
Start with a cost function $J(\beta)$:



Then gradually move towards the minimum.

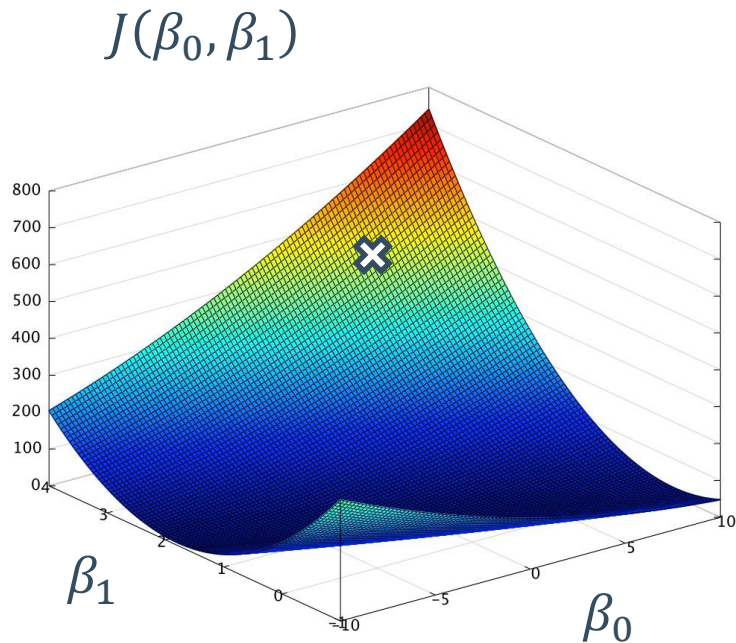
GRADIENT DESCENT WITH LINEAR REGRESSION

- Now imagine there are two parameters (β_0, β_1)
- This is a more complicated surface on which the minimum must be found
- How can we do this without knowing what $J(\beta_0, \beta_1)$ looks like?



GRADIENT DESCENT WITH LINEAR REGRESSION

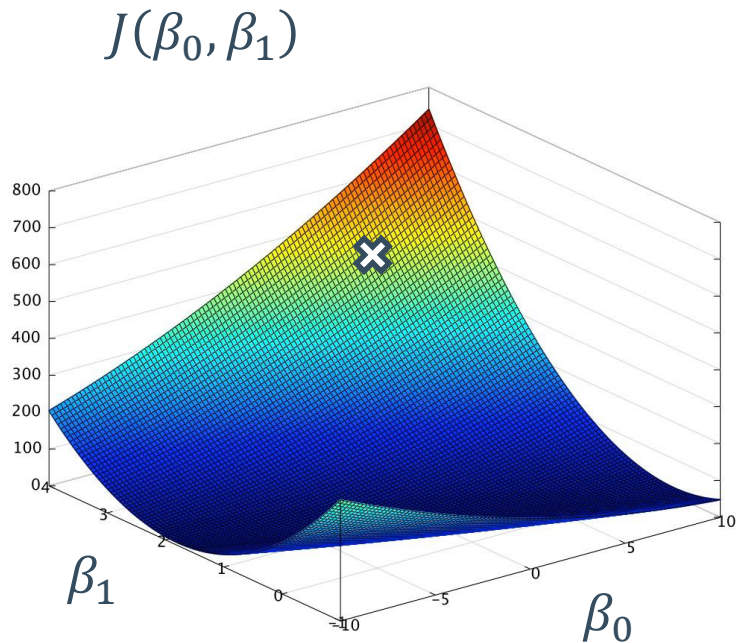
- Compute the gradient, $\nabla J(\beta_0, \beta_1)$, which points in the direction of the biggest increase!
- $-\nabla J(\beta_0, \beta_1)$ (negative gradient) points to the biggest decrease at that point!



GRADIENT DESCENT WITH LINEAR REGRESSION

- The gradient is the a vector whose coordinates consist of the partial derivatives of the parameters

$$\nabla J(\beta_0, \dots, \beta_n) = \left\langle \frac{\partial J}{\partial \beta_0}, \dots, \frac{\partial J}{\partial \beta_n} \right\rangle$$

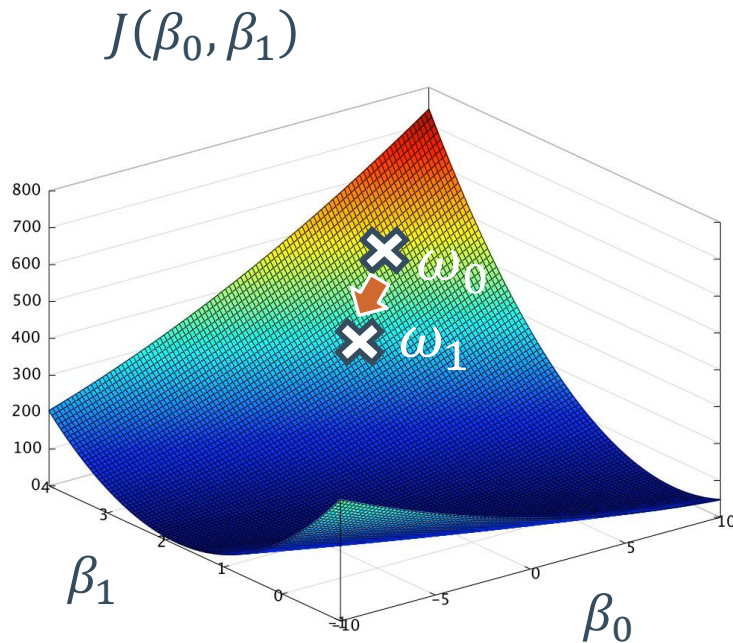


GRADIENT DESCENT WITH LINEAR REGRESSION

- Then use the gradient (∇) and the cost function to calculate the next point (ω_1) from the current one (ω_0):

$$\omega_1 = \omega_0 - \alpha \nabla \frac{1}{2} \sum_{i=1}^m \left((\beta_0 + \beta_1 x_{obs}^{(i)}) - y_{obs}^{(i)} \right)^2$$

- The learning rate (α) is a tunable parameter that determines step size

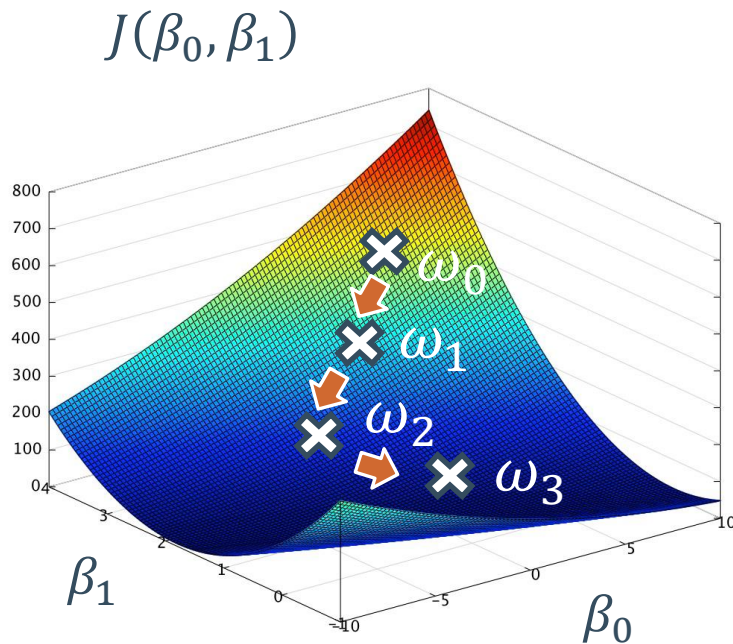


GRADIENT DESCENT WITH LINEAR REGRESSION

- Each point can be iteratively calculated from the previous one

$$\omega_2 = \omega_1 - \alpha \nabla \frac{1}{2} \sum_{i=1}^m \left((\beta_0 + \beta_1 x_{obs}^{(i)}) - y_{obs}^{(i)} \right)^2$$

$$\omega_3 = \omega_2 - \alpha \nabla \frac{1}{2} \sum_{i=1}^m \left((\beta_0 + \beta_1 x_{obs}^{(i)}) - y_{obs}^{(i)} \right)^2$$



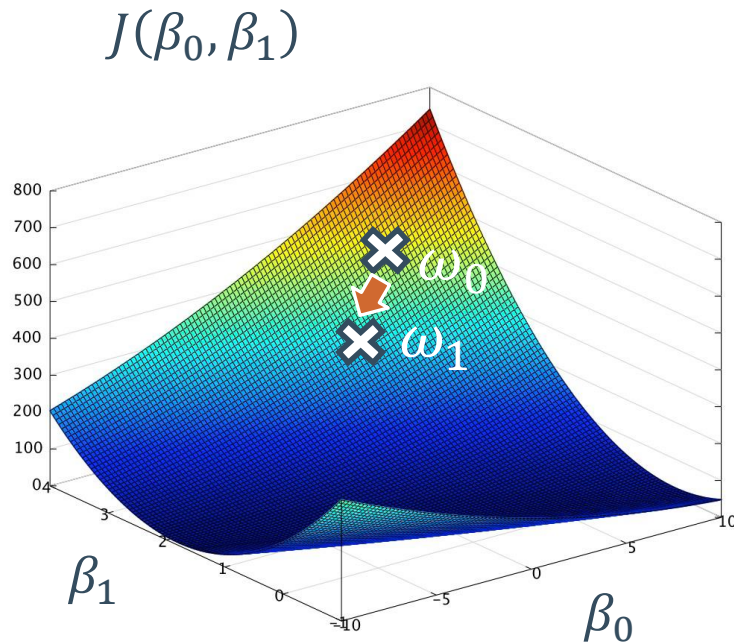
STOCHASTIC GRADIENT DESCENT

- Use a single data point to determine the gradient and cost function instead of all the data

$$\omega_1 = \omega_0 - \alpha \nabla \frac{1}{2} \sum_{i=1}^m \left((\beta_0 + \beta_1 x_{obs}^{(i)}) - y_{obs}^{(i)} \right)^2$$



$$\omega_1 = \omega_0 - \alpha \nabla \frac{1}{2} \left((\beta_0 + \beta_1 x_{obs}^{(0)}) - y_{obs}^{(0)} \right)^2$$

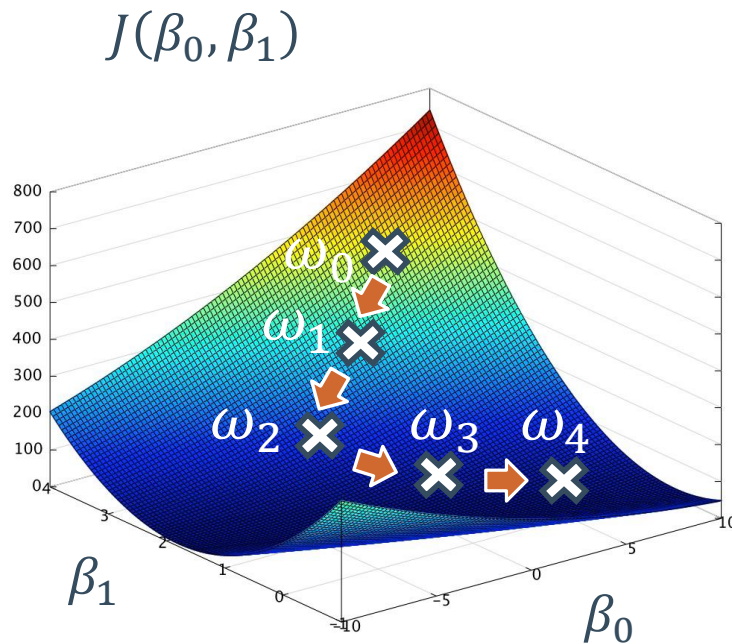


STOCHASTIC GRADIENT DESCENT

- Use a single data point to determine the gradient and cost function instead of all the data

$$\begin{aligned}\omega_1 &= \omega_0 - \alpha \nabla \frac{1}{2} \left((\beta_0 + \beta_1 x_{obs}^{(0)}) - y_{obs}^{(0)} \right)^2 \\ &\dots \\ \omega_4 &= \omega_3 - \alpha \nabla \frac{1}{2} \left((\beta_0 + \beta_1 x_{obs}^{(3)}) - y_{obs}^{(3)} \right)^2\end{aligned}$$

- Path is less direct due to noise in single data point—"stochastic"



MINI BATCH GRADIENT DESCENT

- Perform an update for every n training examples

$$\omega_1 = \omega_0 - \alpha \nabla \frac{1}{2} \sum_{i=1}^n \left((\beta_0 + \beta_1 x_{obs}^{(i)}) - y_{obs}^{(i)} \right)^2$$

Best of both worlds:

- Reduced memory relative to "vanilla" gradient descent
- Less noisy than stochastic gradient descent

