

# CSC 461: Machine Learning

## Fall 2024

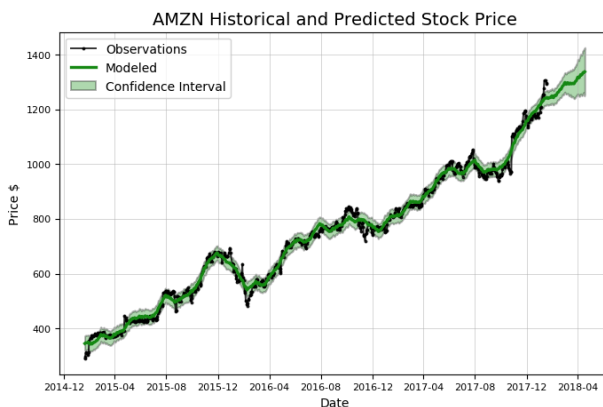
## Linear regression

Prof. Marco Alvarez, Computer Science  
University of Rhode Island

## Preliminaries

## Continuous output

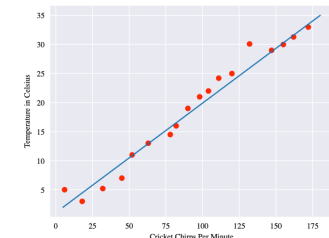
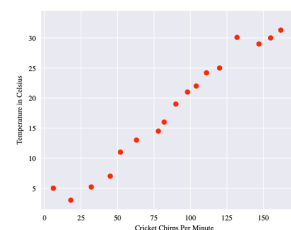
- Certain applications require the prediction of continuous values



<https://towardsdatascience.com/stock-prediction-in-python-b66555171a2>

## Linear functions

- Assumes the output  $y$  is a linear function of the input  $x$
- can use the function to make predictions, very simple approach, e.g. linear regression



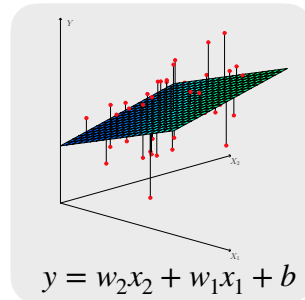
$$y = wx + b \quad w, x, b \in \mathbb{R}$$

slope      intercept

## Linear functions

### What if we have $d$ features?

- $\mathbf{w}, \mathbf{x} \in \mathbb{R}^d$
- $b \in \mathbb{R}$



$$h(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$$

hypothesis      weights      bias

The weights and bias are the model **parameters** which define the hypothesis and are used to make predictions

## Linear regression

## Alternative notation

### Augmented vectors

- incorporate bias into the weight vector
- $\mathbf{w} = [b, w_1, w_2, \dots, w_d]^T$
- augment input vector with 1
- $\mathbf{x} = [1, x_1, x_2, \dots, x_d]^T$

$x_1$	$x_2$	$y$
0.5	0.1	0.25
0.3	0.9	0.5
0.3	0.875	1.15
0.45	0.15	2.13
...	...	...



$x_0$	$x_1$	$x_2$	$y$
1	0.5	0.1	0.25
1	0.3	0.9	0.5
1	0.3	0.875	1.15
1	0.45	0.15	2.13
...	...	...	...

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

Hypothesis space  $\mathcal{H} = \{h_{\mathbf{w}} : h_{\mathbf{w}}(\mathbf{x}) = \mathbf{w}^T \mathbf{x}, \mathbf{w} \in \mathbb{R}^{d+1}\}$

## Linear regression

### Fundamental supervised learning algorithm

- used for predicting continuous target variables
- assumes a linear relationship between input features and the target

### Applications

- financial forecasting (e.g., stock prices, economic indicators)
- environmental science (e.g., climate change predictions)
- marketing (e.g., sales forecasting)
- real estate (e.g., house price prediction)

# Linear regression

## ► Data

- $\mathcal{D} = \{(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(n)}, y^{(n)})\}$
- $\mathbf{x}^{(i)} \in \mathbb{R}^{d+1}, y^{(i)} \in \mathbb{R}$

## ► Model

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

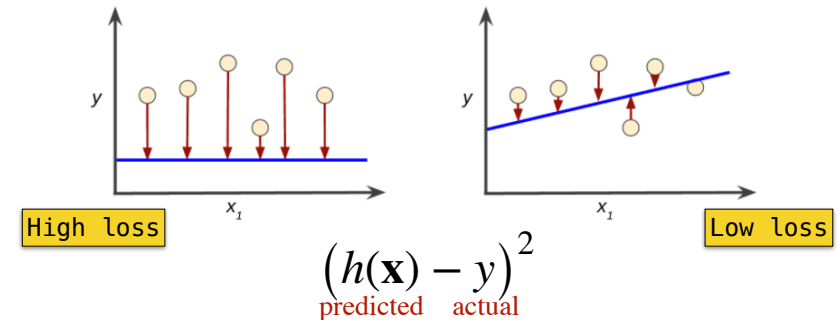
## ► Squared Loss (L2)

- $l_{sq}(h, \mathbf{x}, y) = (h(\mathbf{x}) - y)^2$
- $L(h, \mathcal{D}) = \frac{1}{n} \sum_{i=1}^n l_{sq}(h, \mathbf{x}^{(i)}, y^{(i)})$

# Residuals and MSE loss

## ► Residual

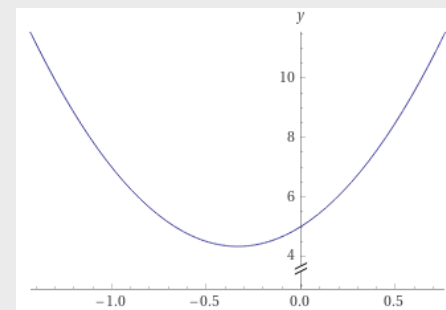
- difference between predicted and actual value



Want to find a **linear function** with **small residuals**

# Closed-form solution

## min vs. argmin



$$f(x) = 6x^2 + 4x + 5$$

$\min f(x)?$

$\arg \min f(x)?$   
 $x$

Given a function, it represents the input value(s) that produce the function's minimum output value

## Normal equations

- Analytical solution to minimize the loss function

$$\mathbf{w}^* = \arg \min_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n (h(\mathbf{x}^{(i)}) - y^{(i)})^2$$



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

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}^{(i)} - y^{(i)})^2$$

## Using matrix notation

$$\mathbf{y} = \mathbf{X}\mathbf{w} \quad \begin{bmatrix} (\mathbf{x}^{(1)})^T \\ \vdots \\ (\mathbf{x}^{(n)})^T \end{bmatrix} \begin{bmatrix} w_0 \\ \vdots \\ w_d \end{bmatrix} \approx \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(n)} \end{bmatrix}$$

## Using matrix notation

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}^{(i)} - y^{(i)})^2$$

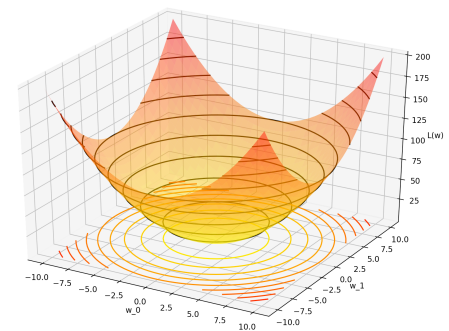


$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \frac{1}{n} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2$$

## Minimizing the loss function

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \frac{1}{n} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2$$

continuous,  
differentiable,  
convex



# Minimizing the loss function

- ▶ Set the gradient to zero and solve for  $\mathbf{w}$

$$L(\mathbf{w}) = \frac{1}{n} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 \quad \nabla L(\mathbf{w}) = \frac{2}{n} \mathbf{X}^T (\mathbf{X}\mathbf{w} - \mathbf{y})$$

$$\frac{2}{n} \mathbf{X}^T (\mathbf{X}\mathbf{w} - \mathbf{y}) = 0$$

$$\mathbf{X}^T \mathbf{X} \mathbf{w} = \mathbf{X}^T \mathbf{y}$$

$$\mathbf{w} = \underbrace{(\mathbf{X}^T \mathbf{X})^{-1}}_{\text{pseudoinverse}} \mathbf{X}^T \mathbf{y}$$

There are other methods for finding the optimal solution  
e.g. gradient descent, MLE

# Show me the code

```
# generate random data
Xtr = np.hstack((np.ones((100, 1)), np.random.rand(100, 4)))
Ytr = np.random.rand(100, 1)
Xte = np.hstack((np.ones((10, 1)), np.random.rand(10, 4)))
Yte = np.random.rand(10, 1)
```

```
w = np.linalg.pinv(Xtr) @ Ytr
pred = Xte @ w
```

or

```
w = np.linalg.inv(Xtr.T @ Xtr) @ Xtr.T @ Ytr
pred = Xte @ w
```

or

```
reg = LinearRegression().fit(Xtr, Ytr)
pred = reg.predict(Xte)
```

```
loss = np.mean((pred-Yte) ** 2)
print(loss)
```

# Conclusion

- ▶ Linear regression: simple yet powerful
  - foundation for many advanced ML techniques
  - serves as a baseline for more complex models
- ▶ Computational Complexity
  - becomes inefficient for large datasets or high-dimensional data
  - time complexity:  $O(nd^2 + d^3)$
- ▶ Alternatives
  - mini-batch gradient descent
  - regularization techniques (Ridge, Lasso)
- ▶ Limitations and Considerations
  - assumes linearity between features and target
  - sensitive to outliers