# EEP 596: Adv Intro ML || Lecture 2
## Dr. Karthik Mohan

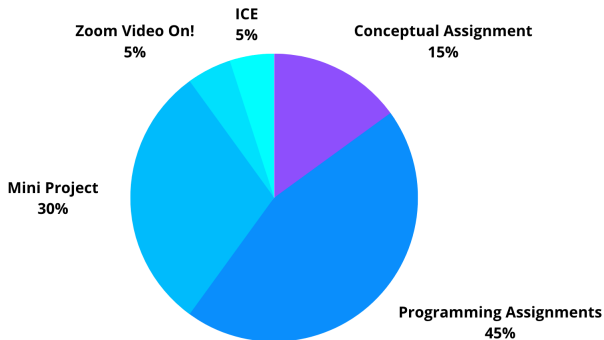Univ. of Washington, Seattle

Jan 6, 2022

# Logistics

- **Lecture** Tuesday Lecture: Expectation that you join in person. Thursday Lecture: Zoom (zoom attendance will be taken).
- **Assignment** Programming Assignment 1 to be assigned - Due next Thursday, January 12th, midnight
- **Office Hours** Karthik: 6 - 6:30 pm on Thursday, Ayush - TBD
- **Calendly slots** Feel free to pick calendly slots for 1:1 15 min syncs as needed (recommended)
- **Course Webpage** https://bytesizeml.github.io/ml2023/

# Weekly Schedule

|  | Day | Timings | Class type |
|---|---|---|---|
| **Lecture 1 (In-person)** | T | 4 pm - 6 pm | In-person |
| **Lecture 2** | Th | 4 pm - 6 pm | Zoom |
| **Office Hours Karthik** | Th | 6 - 6:30 pm | Zoom |
| **Office Hours Ayush** | TBD | TBD | Zoom |
| **Quiz Section Ayush** | TBD | TBD | Zoom |
| **Grading hours** | TBD | TBD | Zoom |

# Assessments Breakdown

# Lecture Structure

Format for each lecture

- Sprinkle in a few In-class exercises MCQ for conceptual understanding
- Where required - Will set extra context on applications/background - This may be slow for some but super useful for rest of class - Let's adjust and adapt!
- Break at 1 hour mark
- Break-outs in between/end of class for peer discussion + networking
- Anything else ?

# Class goes at the average pace!

Quick pointers

- We will cater the lecture to discuss fundamentals and go at a pace comfortable with the average of the class
- If the class/topic is going too fast for you - There maybe brushing up of background (e.g. linear algebra/calculus/programming) that you may have to do in your own time!
- If a topic is going slow - Opportunity to dive deeper into the topic through additional reading of papers or programming
- Be sure to brush up/catch up on your python and linear algebra to gear up for upcoming lectures and assignments

# Lectures and Programming Assignments (Tentatively)

| Week | Lecture Material | Assignment |
|---|---|---|
| 1 | Linear Regression | Housing Price Prediction |
| 2 | Classification | Spam classification (Kaggle) |
| 3 | Classification | Flower/Leaf classification |
| 4 | Clustering | MNIST digits clustering |
| 5 | Anomaly Detection | Crypto Prediction (Kaggle + P) |
| 6 | Data Visualization | Crypto Prediction (Kaggle + P) |
| 7 | Deep Learning | Visualizing 1000 images |
| 8 | Deep Learning (DL) | ECG Arrythmia Detection |
| 9 | DL in NLP | TwitterSentiment Analysis (Kaggle + P) |
| 10 | DLs in Vision | TwitterSentiment Analysis (Kaggle + P) |

# Coding pointers

- Assignments assume python as the main language (e.g. for hints and modules, etc)
- Coding environment set-up will be one of the problems on HW 1
- Prototyping can be done on notebooks and submitted as such for smaller assignments.
- For mini-projects and kaggle assignments - Please keep your code modular and organized.

# Coding Environment

- Pointers below if you want to get set up on Google Colab for both prototyping, running machine-intensive ML experiments and working with code through IDEs
- Prototype Coding work in Notebooks recommended on Google Colab
- For terminal access on Google Colab, sign up for pro
- `pip3 install colabcode` on termainal
- ColabCode enables you to have a VSCode IDE port into Google Colab - So you can work on the IDE from your laptop but run experiments on Google Colab!

# Maximizing Your Learning of Machine Learning!

- Ask questions during lectures - Clarify things as they happen!

# Maximizing Your Learning of Machine Learning!

- Ask questions during lectures - Clarify things as they happen!
- Make use of office hours and quiz section!

# Maximizing Your Learning of Machine Learning!

- Ask questions during lectures - Clarify things as they happen!
- Make use of office hours and quiz section!
- Collaborative learning - Discord is a great place to brainstorm concepts outside class and unblock yourself.

# Maximizing Your Learning of Machine Learning!

- Ask questions during lectures - Clarify things as they happen!
- Make use of office hours and quiz section!
- Collaborative learning - Discord is a great place to brainstorm concepts outside class and unblock yourself.
- 30% of your learning happens in class and office hours - The remaining 70% happen when you work on the assignments. (You ofcourse need the 30 to get to the 70 :D)
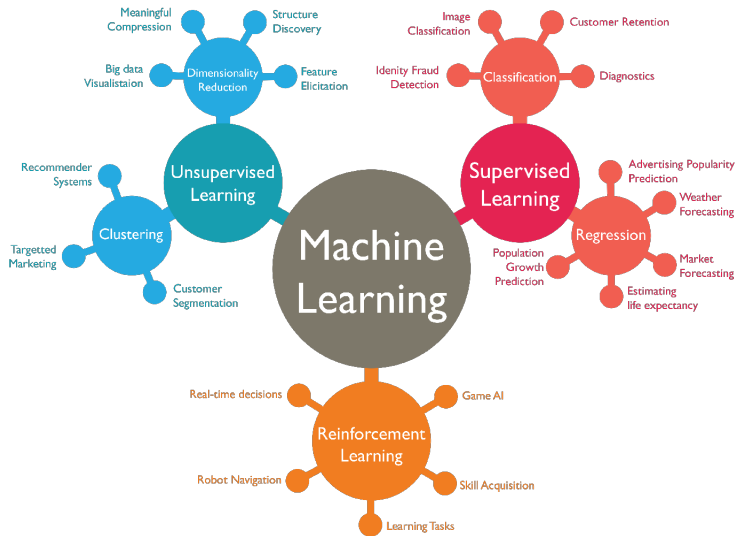
# Maximizing Your Learning of Machine Learning!

- Ask questions during lectures - Clarify things as they happen!
- Make use of office hours and quiz section!
- Collaborative learning - Discord is a great place to brainstorm concepts outside class and unblock yourself.
- 30% of your learning happens in class and office hours - The remaining 70% happen when you work on the assignments. (You ofcourse need the 30 to get to the 70 :D)
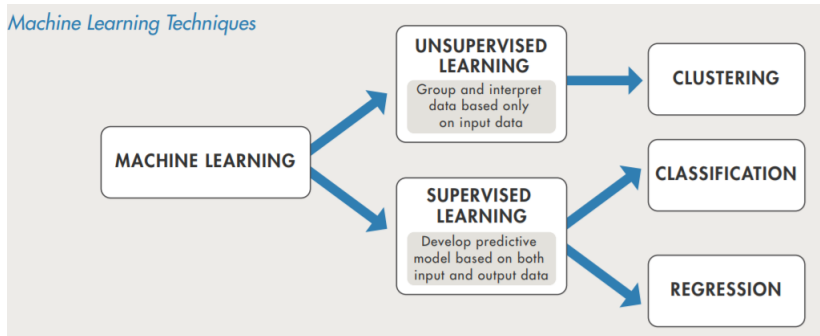- What you put in is what you get out!

# Maximizing Your Learning of Machine Learning!

- Ask questions during lectures - Clarify things as they happen!
- Make use of office hours and quiz section!
- Collaborative learning - Discord is a great place to brainstorm concepts outside class and unblock yourself.
- 30% of your learning happens in class and office hours - The remaining 70% happen when you work on the assignments. (You ofcourse need the 30 to get to the 70 :D)
- What you put in is what you get out!
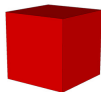- Excitement + Smart work + Inquisitiveness = Maximized learning!
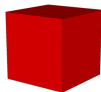
# What is Machine Learning?

# Supervised vs Unsupervised Learning

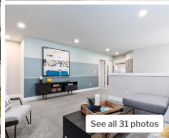# Supervised Learning

# Un-Supervised Learning

# Our first ML method: Linear Regression

# Application: Housing Prices

# Redfin

# Redfin Estimate

**This home is popular**
It's been viewed 2,022 times. Tour it in person or via video chat before it's gone!

Today: 6:00 pm · 7:00 pm · 8:00 pm · More times

## Go tour this home

| TUESDAY 3 JAN | WEDNESDAY 4 JAN | THURSDAY 5 JAN |
|---|---|---|

Tour in person · Tour via video chat

**Schedule tour**

It's free, with no obligation — cancel anytime.

OR

**Start an offer**

Ask a question · (425) 584-3263

## About This Home

Pacific Ridge presents Ironwood! Gorgeous new home community centrally located between Bothell, Mill Creek & Lynnwood. Perched just off North Road with panoramic views to the East, this neighborhood offers a quiet place to call home with community parks & convenient access to

Continue reading ⌄

Listed by Mari Riksheim · Pacific Ridge - DRH, LLC
Listed by Melissa Cogswell · Pacific Ridge - DRH, LLC
Redfin checked: 3 minutes ago (Jan 3, 2023 at 2:57pm) · Source: NWMLS #2024145

### Home Facts

| | | | |
|---|---|---|---|
| Status | Active | Time on Redfin | 5 days |
| Property Type | Residential, Residential | HOA Dues | $88/month |
| Year Built | 2023 | Style | Contemporary |
| Community | Lynnwood | Lot Size | 6,252 Sq. Ft. |
| MLS# | 2024145 | | |

### Price Insights

| | | | |
|---|---|---|---|
| List Price | $1,134,995 | Est. Mo. Payment | $7,420 |
| Redfin Estimate | $1,136,063 | Price/Sq.Ft. | $420 |

# Zoom Breakout #1

### Zillow Estimate/RedFin Estimate

If you are on the market to buy a house, you would perhaps be looking at "Zestimates" or "RedFin Estimates" to filter out houses in your budget range. Discuss in your group, what are the factors that influence the price of a home and what are the factors (also called features in ML) that may have been used to construct these estimates. Once you have a set of factors identified, how do you combine them to produce the final house price estimate?

# Typical Housing Data, Seattle

| Index | SqFt | #Rooms | # Bathrooms | Location | Selling Price |
|-------|------|--------|-------------|----------|---------------|
| 1 | 2500 | 4 | 3 | Bothell | 1 MM |
| 2 | 2000 | 3 | 2 | Bellevue | 950k |
| 3 | 3000 | 4 | 3 | Sammamish | 1.3 MM |
| 4 | 3000 | 4 | 3 | Issaquah High | 1.6 MM |
| 5 | . . . | . . . | . . . | . . . | . . . |
| | | | | | |

# Typical Housing Data, Seattle

| Index | SqFt | #Rooms | # Bathrooms | Location | **Selling Price** |
|-------|------|--------|-------------|----------|-------------------|
| 1 | 2500 | 4 | 3 | Bothell | 1 MM |
| 2 | 2000 | 3 | 2 | Bellevue | 950k |
| 3 | 3000 | 4 | 3 | Sammamish | 1.3 MM |
| 4 | 3000 | 4 | 3 | Issaquah High | 1.6 MM |
| 5 | . . . | . . . | . . . | . . . | . . . |
| | | | | | |

# Other attributes for housing price prediction

Other attributes that matter?

# Categorical vs Numerical Attributes

## Categorical

Attributes that fall into a clear set of categories. Example: zipcode of a place

# Categorical vs Numerical Attributes

### Categorical
Attributes that fall into a clear set of categories. Example: zipcode of a place

### Numerical
Attributes that fall in a numeric range. Example: weight or height of a person

# Categorical vs Numerical Attributes

## Categorical

Attributes that fall into a clear set of categories. Example: zipcode of a place

## Numerical

Attributes that fall in a numeric range. Example: weight or height of a person

## Modeling Choice

Sometimes, whether an attribute is categorical or numerical is a modeling choice!

# Categorical vs Numerical Attributes

Categorical or Numerical??

| Index | SqFt | #Rooms | # Bathrooms | Location | Selling Price |
|-------|------|--------|-------------|----------|---------------|
| 1 | 2500 | 4 | 3 | Bothell | 1 MM |
| 2 | 2000 | 3 | 2 | Bellevue | 950k |
| 3 | 3000 | 4 | 3 | Sammamish | 1.3 MM |
| 4 | 3000 | 4 | 3 | Issaquah High | 1.6 MM |
| 5 | . . . | . . . | . . . | . . . | . . . |
| | | | | | |

# Matrices and Vectors

### Data matrix $X$

Let's say in our Housing database, we have 1000 houses and 30 attributes. If we wanted to represent this as a data matrix, $X$, what would be the dimensions of such a matrix ?

# Matrices and Vectors

## Data matrix $X$

Let's say in our Housing database, we have 1000 houses and 30 attributes. If we wanted to represent this as a data matrix, $X$, what would be the dimensions of such a matrix ?

## Price vector $y$

For the same example as before, we take the housing prices of all the homes and put them into a price vector $y$. What would be the dimension of this vector $y$ ?

# $X$ and $y$ in housing data

| Index | SqFt | #Rooms | # Bathrooms | Location | **Selling Price** |
|-------|------|--------|-------------|----------|-------------------|
| 1 | 2500 | 4 | 3 | Bothell | 1 MM |
| 2 | 2000 | 3 | 2 | Bellevue | 950k |
| 3 | 3000 | 4 | 3 | Sammamish | 1.3 MM |
| 4 | 3000 | 4 | 3 | Issaquah High | 1.6 MM |
| 5 | . . . | . . . | . . . | . . . | . . . |
| | | | | | |

# Linear Model

## Linear Model

In Linear models, we assume that the target, $y$ is a linear combination of the attributes or features $x$. This is a 'modeling assumption'. The combination is represented by a weight vector $w$.
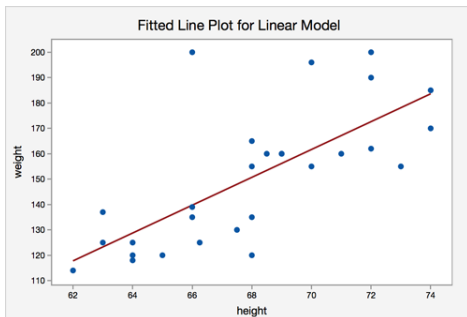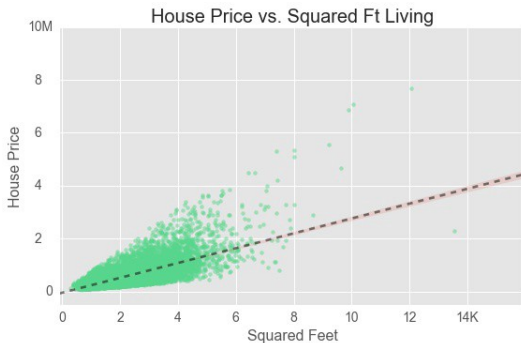
# Linear Model

## Linear Model

In Linear models, we assume that the target, $y$ is a linear combination of the attributes or features $x$. This is a 'modeling assumption'. The combination is represented by a weight vector $w$.

## Visualizing Linear Model

# Linear Model for Housing Prices Application



House Price vs. Squared Ft Living

## Linear Model

In Linear models, we assume that the target, $y$ is a linear combination of the attributes or features $x$. The combination is represented by a weight vector $w$.

## Linear Model

In Linear models, we assume that the target, $y$ is a linear combination of the attributes or features $x$. The combination is represented by a weight vector $w$.

## In the housing price example

$$y = w_0 + w_1 \times x_1 + w_2 \times x_2 + w_3 \times x_3 + w_4 \times x_4$$

## Linear Model

In Linear models, we assume that the target, $y$ is a linear combination of the attributes or features $x$. The combination is represented by a weight vector $w$.

### In the housing price example

$$y = w_0 + w_1 \times x_1 + w_2 \times x_2 + w_3 \times x_3 + w_4 \times x_4$$

### In the housing price example

$$1MM = w_0 + w_1 \times 2500 + w_2 \times 4 + w_3 \times 3 + w_4 \times \text{ Bothell}$$

### Linear Model

In Linear models, we assume that the target, $y$ is a linear combination of the attributes or features $x$. The combination is represented by a weight vector $w$.

### In the housing price example

$$y = w_0 + w_1 \times x_1 + w_2 \times x_2 + w_3 \times x_3 + w_4 \times x_4$$

### In the housing price example

$$1MM = w_0 + w_1 \times 2500 + w_2 \times 4 + w_3 \times 3 + w_4 \times \text{Bothell}$$

### There's one problem though!

How do we multiply a 'location' by a weight ?

# Dealing with categorical attributes

**One approach: Create new dummy attributes!**

$x_{Bothell}, x_{Bellevue}, x_{Sammamish}, x_{IssaquahHigh}$ - One dummy variable for each location that takes a value 1 if its the true location and 0 otherwise.

# Dealing with categorical attributes

## One approach: Create new dummy attributes!

$x_{Bothell}, x_{Bellevue}, x_{Sammamish}, x_{IssaquahHigh}$ - One dummy variable for each location that takes a value 1 if its the true location and 0 otherwise.

## ICE #1 (2 mins): How many attribues do we have now?

Let's say our data consisted of the following attributes: Square Footage, # Rooms, # Bathrooms, Location. After applying "pre-processing" to the data of introducing dummy attributes, how many total attributes do we have now ? Answer poll pollev.com/karthikmohan088

# Modifying the Data Matrix

Where we started: $X$

| Index | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $y$ |
|-------|-------|-------|-------|---------|------|
| 1 | 2500 | 4 | 3 | Bothell | 1 MM |

# Modifying the Data Matrix

Where we started: $X$

| Index | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $y$ |
|-------|-------|-------|-------|---------|------|
| 1 | 2500 | 4 | 3 | Bothell | 1 MM |

After pre-processing for categorical attributes: New $X$

| Index | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ | $x_6$ | |
|-------|-------|-------|-------|-------|-------|-------|--|
| 1 | | | | | | | |
| 2 | | | | | | | |
| ⋮ | | | | | | | |

# Modifying the Data Matrix

Where we started: $X$

| Index | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $y$ |
|-------|-------|-------|-------|---------|------|
| 1 | 2500 | 4 | 3 | Bothell | 1 MM |

After pre-processing for categorical attributes: New $X$

| Index | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ | $x_6$ | |
|-------|-------|-------|-------|-------|-------|-------|---|
| 1 | | | | | | | |
| 2 | | | | | | | |
| ⋮ | | | | | | | |

Does vector $y$ change?

# Back to the Linear Model

## A formula for the house price

Let $y_i$ be the price of the $i_{th}$ home. Let $X_{ij}$ denote the $j_{th}$ attribute of the $i_{th}$ home. Then

$$y_i \sim w_0 + w_1 \times X_{i1} + w_2 \times X_{i2} + w_3 \times X_{i3} + \ldots$$

# Back to the Linear Model

### A formula for the house price

Let $y_i$ be the price of the $i_{th}$ home. Let $X_{ij}$ denote the $j_{th}$ attribute of the $i_{th}$ home. Then

$$y_i \sim w_0 + w_1 \times X_{i1} + w_2 \times X_{i2} + w_3 \times X_{i3} + \ldots$$

### A succinct expression for the $i_{th}$ house

$$y_i = w^T X_{i,.} = w \cdot X_{i,.}$$

# Back to the Linear Model

A formula for the house price

Let $y_i$ be the price of the $i_{th}$ home. Let $X_{ij}$ denote the $j_{th}$ attribute of the $i_{th}$ home. Then

$$y_i \sim w_0 + w_1 \times X_{i1} + w_2 \times X_{i2} + w_3 \times X_{i3} + \ldots$$

A succinct expression for the $i_{th}$ house

$$y_i = w^T X_{i,.} = w \cdot X_{i,.}$$

ICE #2 (2 mins): Succinct expression for $y$ in terms of $X$ and $w$?

# Linear Regression: Putting it all Together

## Definition

Find the best weights/parameters/coefficients $w$ such that $X_{i,:}^T w$ is as close to $y_i$ as possible!

# Linear Regression: Putting it all Together

### Definition

Find the best weights/parameters/coefficients $w$ such that $X_{i,:}^T w$ is as close to $y_i$ as possible!

### Mathematically

Minimize the following expression:

$$\min_w \|Xw - y\|_2^2$$

# Linear Regression: Putting it all Together

### Definition

Find the best weights/parameters/coefficients $w$ such that $X_{i,.}^T w$ is as close to $y_i$ as possible!

### Mathematically

Minimize the following expression:

$$\min_w \|Xw - y\|_2^2$$

### Estimate or "learned" parameter

Represented usually by $\hat{w}$ and $\hat{y}$ is the "predicted" house price for all the homes.

# Linear Regression: Putting it all Together

## Definition

Find the best weights/parameters/coefficients $w$ such that $X_{i,:}^T w$ is as close to $y_i$ as possible!

## Mathematically

Minimize the following expression:

$$\min_w \|Xw - y\|_2^2$$

## Estimate or "learned" parameter

Represented usually by $\hat{w}$ and $\hat{y}$ is the "predicted" house price for all the homes.

## ICE #3 (1 min)

What's the succinct expression for $\hat{y}$?
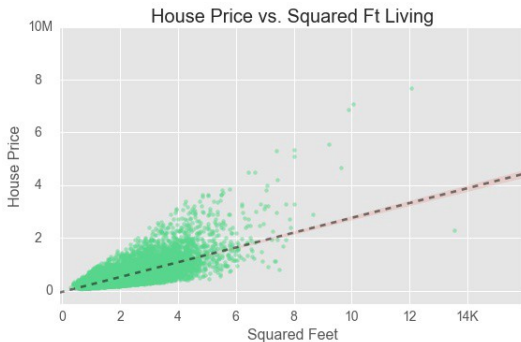
# Line of best fit

## Best fit

$\hat{w}$ defines the line of best fit. $h(x) = \hat{w}^T x$ gives us the line and in higher dimensions, it's called a "hyperplane".

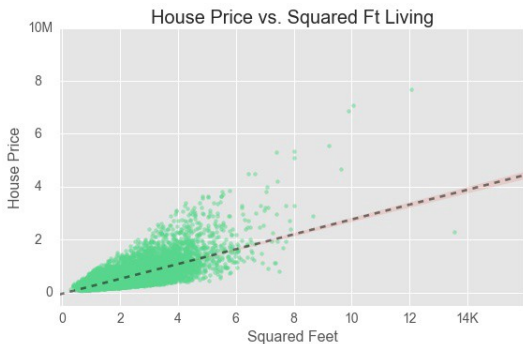# Line of best fit

## Best fit

$\hat{w}$ defines the line of best fit. $h(x) = \hat{w}^T x$ gives us the line and in higher dimensions, it's called a "hyperplane".

## Housing price example

# Line of best fit

## Housing price example
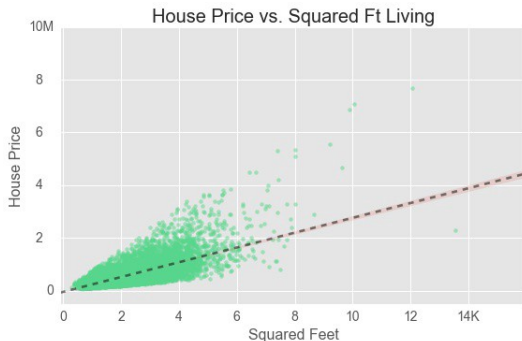


House Price vs. Squared Ft Living

## ICE #4 (1 min)

What would you say is the value of the bias, $w_0$ for the line in the visual above?

# Hyperplane in 3 dimensions

## Housing price example



## 3 dim hyperplane

$$\hat{y} = w_0 + w_1 \times x_1 + w_2 \times x_2$$

$x_1$ could be square footage and $x_2$ could be number of bedrooms.

# Linear Regression

## Closed form!

There is actually a closed form expression for Linear Regression!

$$\min_w \|Xw - y\|_2^2$$

$\hat{w} = (X^T X)^{-1} X^T y$! (Q: How do we arrive at this?)

# Linear Regression

## Closed form!

There is actually a closed form expression for Linear Regression!

$$\min_w \|Xw - y\|_2^2$$

$\hat{w} = (X^T X)^{-1} X^T y$! (Q: How do we arrive at this?)

## In practice!

In practice, a linear regression library might revert to doing "gradient descent" on the learning objective. Why do that?

# Housing price Example

## Pre-processing of data

One is taking care of categorical variables such as location with dummy attributes (also called 'bag of words' model). Anything else we may need to do on the data to get good predictions?

# Training the Linear Regression Model

| $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ | $x_6$ | $x_7$ | $y$ |
|---|---|---|---|---|---|---|---|
| | | | | | | | |
| | | | | | | | |
| | | | | | | | |
| | | | | | | | |
| | | | | | | | |
| | | | | | | | |
| | | | | | | | |
| | | | | | | | |
| | | | | | | | |
| | | | | | | | |

# Can we use of all of data for training?

- Why not use all data for training ?

# The phenomenon of Overfitting

### Overfitting

Overfitting is when your mdoel performs great on training data but doesn't match up on test data. To account for overfitting, we also have a validation data set.

# Understanding over-fitting better

When do we expect over-fitting?

When the number of attributes in our model exceeds the size of the data set.

In terms of data matrix $X$

# rows $<<$ # columns

- **Training** Learning parameters/weights, i.e. $w$ for Linear Regression is called Training.

# Data Splits for Machine Learning

- **Training** Learning parameters/weights, i.e. $w$ for Linear Regression is called Training.
- We don't use all data for training - Some portion of the data is kept for validation and testing.

# Data Splits for Machine Learning

- **Training** Learning parameters/weights, i.e. $w$ for Linear Regression is called Training.
- We don't use all data for training - Some portion of the data is kept for validation and testing.
- **Data Splits:** Usually, 80% of data is kept for training, 10% for validation and 10% for training. The splits are chosen randomly.

- **Training** Learning parameters/weights, i.e. $w$ for Linear Regression is called Training.
- We don't use all data for training - Some portion of the data is kept for validation and testing.
- **Data Splits:** Usually, 80% of data is kept for training, 10% for validation and 10% for training. The splits are chosen randomly.
- Why not use all data for training ?

# Data Splits for Machine Learning

- **Training** Learning parameters/weights, i.e. $w$ for Linear Regression is called Training.
- We don't use all data for training - Some portion of the data is kept for validation and testing.
- **Data Splits:** Usually, 80% of data is kept for training, 10% for validation and 10% for training. The splits are chosen randomly.
- Why not use all data for training ?
- Why not just have **train** and **test** data? What's the point of validation data set?

# Example: 70 : 10 : 20 Train-Val-Test data split

Choose 70% train data at random

| $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ | $x_6$ | $x_7$ | $y$ |
|---|---|---|---|---|---|---|---|
| | | | | | | | |
| | | | | | | | |
| | | | | | | | |
| | | | | | | | |
| | | | | | | | |
| | | | | | | | |
| | | | | | | | |
| | | | | | | | |
| | | | | | | | |
| | | | | | | | |

# Example: 70 : 10 : 20 Train-Val-Test data split

Add 20% test data at random

# Example: 70 : 10 : 20 Train-Val-Test data split

Remainder becomes validation data

| $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ | $x_6$ | $x_7$ | $y$ |
|---|---|---|---|---|---|---|---|
| | | | | | | | |
| | | | | | | | |
| | | | | | | | |
| | | | | | | | |
| | | | | | | | |
| | | | | | | | |
| | | | | | | | |
| | | | | | | | |
| | | | | | | | |
| | | | | | | | |

# Coding Pointers

- **Pandas** library in Python is good for data pre-processing before training your Linear Regression model

# Coding Pointers

- **Pandas** library in Python is good for data pre-processing before training your Linear Regression model
- **Dummy attributes** for categorical variables can also be added in through `pandas.get_dummies()` method.

# Coding Pointers

- **Pandas** library in Python is good for data pre-processing before training your Linear Regression model
- **Dummy attributes** for categorical variables can also be added in through `pandas.get_dummies()` method.
- Use **Scikit-learn** for implementing Linear Regression

# Coding Pointers

- **Pandas** library in Python is good for data pre-processing before training your Linear Regression model
- **Dummy attributes** for categorical variables can also be added in through `pandas.get_dummies()` method.
- Use **Scikit-learn** for implementing Linear Regression
- Should now be ready to tackle both the conceptual and programming Assignment 1!

# Summary so far

- Linear Regression finds a line of best fit through the data.
- $R^2$ measure determines the goodness of fit.
- Usually multiple good attributes are needed for a good prediction and a good fit.
- Data pre-processing. Categorical attributes are handled through creation of dummy attributes and in addition normalizing of the attributes brings all attributes on the same scale for regression.
- We have a closed form/analytical solution for Linear Regression, but for large data sets, gradient descent algorithm (iterative) gets used for scalability reasons.
- We don't use all of a data set for training. A portion of data is kept for validation and testing. This is to prevent over-fitting and also for fair evaluation purposes.
- The data set split is usually $80 - 10 - 10$ or $70 - 10 - 20$ (train-val-test).

# Summary so far

- Over-fitting happens when we have fewer data points as compared to the number of attributes or features.
- Over-fitting can be taken care off by increasing data-set size, decreasing number of attributes or through regularization strategies

# Questions/Thoughts?