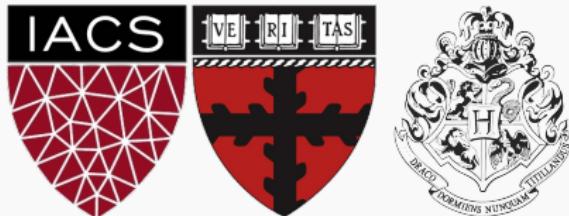


Workshop #1: Introduction to Python for Data Science

TRiCAM 2017

W. Pan



Lecture Outline

Review

Motivating Regression

Stats Review

Linear Regression (Univariate)

Linear Regression (Multivariate)

Polynomial Regression

What is Machine Learning?

Review

The Data Science Process

Recall the Data Science Process we outlined yesterday:

- ▶ Ask questions
- ▶ Data Collection
- ▶ Data Exploration
- ▶ **Data Modeling**
- ▶ **Data Analysis**
- ▶ Visualization and Presentation of Results

Yesterday we addressed data collection and exploration. Today we'll be addressing building models for data and analyzing the results.

Data Collection: Summary

A couple of important observations from yesterday:

- ▶ Real world data can be a mess. It can have:
 - missing values (requires data imputation)
 - erroneous values (requires outlier detection)
 - inconsistent or non-comparable units (requires rescaling, normalization or standardization)
 - messy format (require reformatting)
- ▶ Real word data can be hard to get. Limitations include:
 - expensive API
 - no API, requires scraping
 - is not digitized
 - very old, or not old enough

Motto: Data scientists put data first! Be vigilant for any potential weakness in the integrity of your data and be willing to spend substantial creative energy on data gathering/processing alone.

Data Exploration: Summary

A couple of important observations from yesterday:

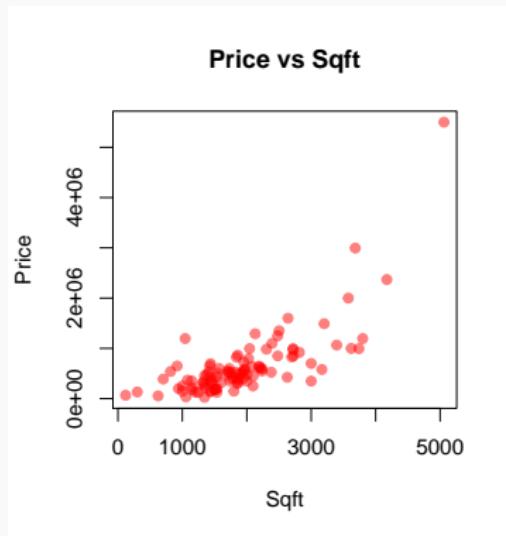
- ▶ **Basic intuition:** shape of data (no. of observations, dimension of feature space), types of variables (quantitative vs categorical).
- ▶ **Descriptive statistics:** summarizing the values
 - **'typical value':** mean or median for quantitative, mode for categorical
 - **'spread':** range, variance/standard deviation, quantiles
- ▶ **Data visualization:** qualitatively exploring relationships and trends
 - **Distribution:** histogram for single variable, scatter plot for multiple variables
 - **Relationships:** scatter plot, trend line
 - **Composition:** pie chart, bar chart, stacked area graph
 - **Comparisons:** multiple trends or histograms in one chart

Try It Yourself!

Apply what you've learned about the data science process to explore a SoCal real-estate listings data set!

Motivating Regression

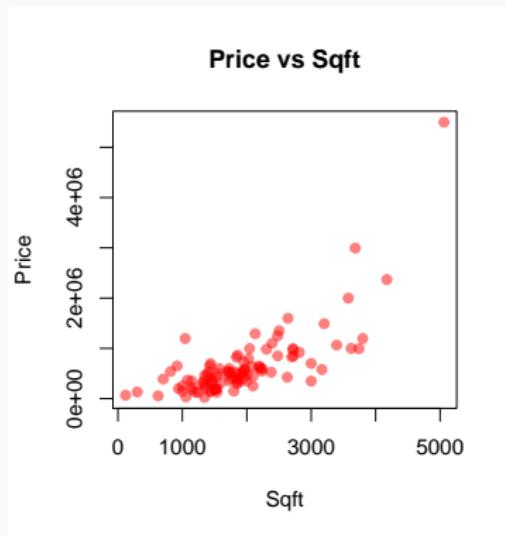
Californian Home Prices (2009)



This is a scatter plot of home prices vs square footage of some homes in southern California.

Can you see any patterns or trends?

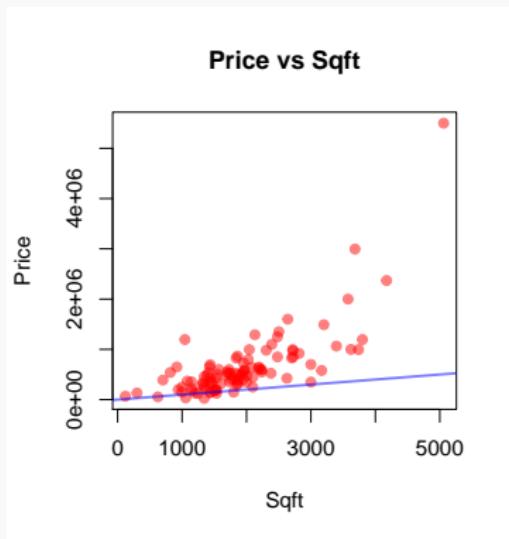
Californian Home Prices (2009)



We see that **as square footage increases, so does price.**

But what is a precise, mathematical description of this relationship?

Californian Home Prices (2009)

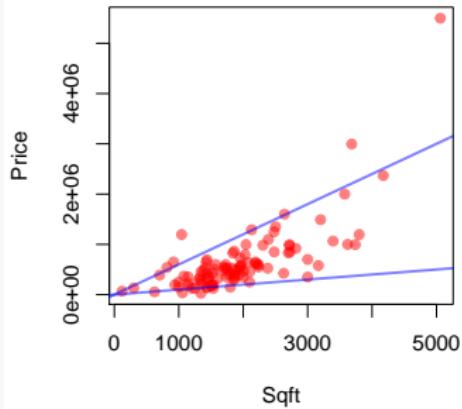


Maybe we want to model the relationship between square footage and price using a simple line.

Does this line capture the trend in the data?

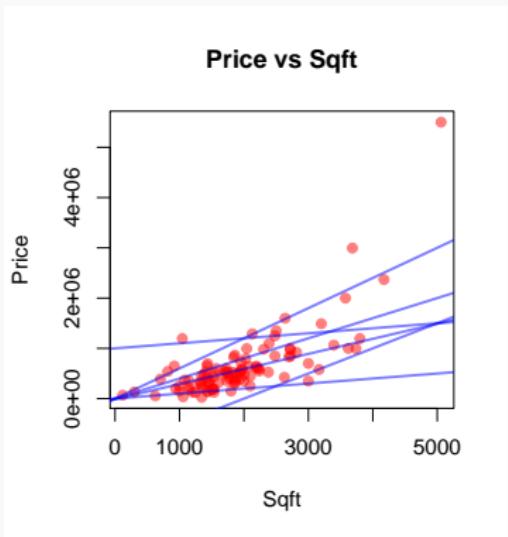
Californian Home Prices (2009)

Price vs Sqft



What about this line?

Californian Home Prices (2009)

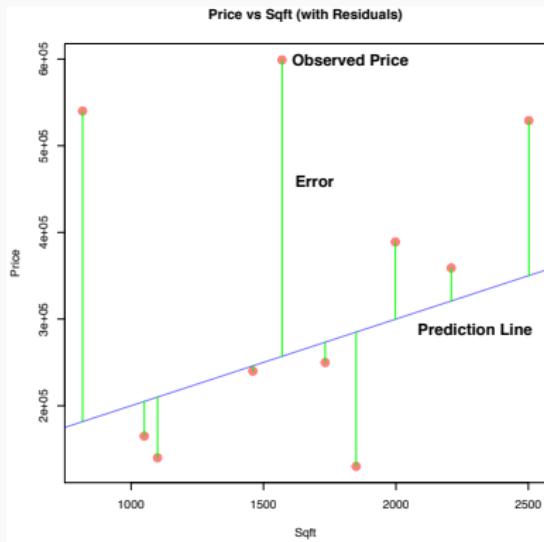


In fact, there are infinite number of lines we can draw through the data.

Which is the best line?

What is a good definition for '**the best line**'?

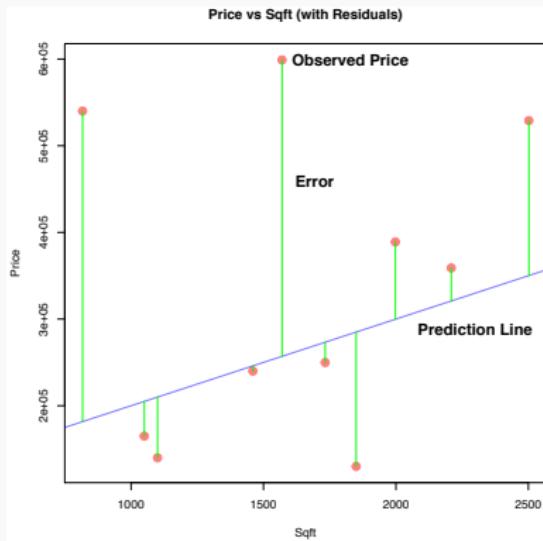
Notion of Error



An **absolute residual** is the absolute difference between the actual price of a home and the price predicted by the line for a given square footage.

$$\text{Res}_i = |\text{Observed}_i - \text{Predicted}_i|$$

Notion of Error



The i -th absolute residual measures the magnitude of the 'error' made by the i -th prediction.

Notions of Fitness

Question: How do we quantify the overall error?

1. **(Max absolute deviation)** Count only the biggest 'error'

$$\max_i |\text{Observed}_i - \text{Predicted}_i|$$

2. **(Sum of absolute deviations)** Add up the 'errors'

$$\sum_i |\text{Observed}_i - \text{Predicted}_i|$$

We can also average them.

3. **(Sum of squared errors)** Add up the squares of the 'errors'

$$\sum_i |\text{Observed}_i - \text{Predicted}_i|^2$$

We can also average them.

Model Fitting

Question: What do we mean by choosing ‘the best line’?

Answer: A line which minimizes the overall error.

Example: Given a set of points $(x_1, y_1), \dots, (x_n, y_n)$, the **average of absolute deviations** of a line $y = mx + b$ is

$$L(m, b) = \frac{1}{n} \sum_{i=1}^n |y_i - (mx_i + b)|$$

L is called the **loss or cost function**. Our goal is to **find \hat{m} and \hat{b} such that the loss, $L(\hat{m}, \hat{b})$, is minimal**:

$$(\hat{m}, \hat{b}) = \operatorname{argmin}_{m,b} L(m, b).$$

Finding the optimal values (\hat{m}, \hat{b}) is called **fitting the linear model**.

Model Fitting

Question: What do we mean by choosing ‘the best line’?

Answer: A line which minimizes the overall error.

Example: Given a set of points $(x_1, y_1), \dots, (x_n, y_n)$, the **mean squared error (MSE)** of a line $y = mx + b$ is

$$L(m, b) = \frac{1}{n} \sum_{i=1}^n |y_i - (mx_i + b)|^2$$

L is called the lost (or cost) function. Our goal is to **find \hat{m} and \hat{b} such that the lost, $L(\hat{m}, \hat{b})$, is minimal**:

$$(\hat{m}, \hat{b}) = \operatorname{argmin}_{m,b} L(m, b).$$

Finding the optimal values (\hat{m}, \hat{b}) is called **fitting the linear model**.

Choosing a Fitness Criterion

Question: What do we mean by ‘the best line’?

Answer: A line which minimizes the overall error.

But which notion of error should we choose (max absolute deviation, sum/average of absolute deviation or sum/average of squared errors)?

The answer depends on **how**, we believe, the ‘residual’ (difference between observed and predicted values) arise.

Choosing a Fitness Criterion

Our belief: The relationship between **price** (P) and **square footage** (A) is linear

$$P = m \cdot A + b \quad (\text{model for theoretical prices})$$

But, in real-life, due to unpredictable circumstances observed prices differ from our pricing rule by some **random** amount, ϵ . This random deviation is called **noise**. So our model for observed housing prices is

$$P = m \cdot A + b + \epsilon \quad (\text{model for observed prices})$$

A model that accounts for uncertainty or randomness, where the output (P) is not deterministically dependent on input (A), is called a **statistical model**. The noise, ϵ , is a **random variable**.

Stats Review

Random Variables

A **random variable** (RV) is a variable whose possible values are numerical outcomes of a random phenomenon. There are two types of random variables:

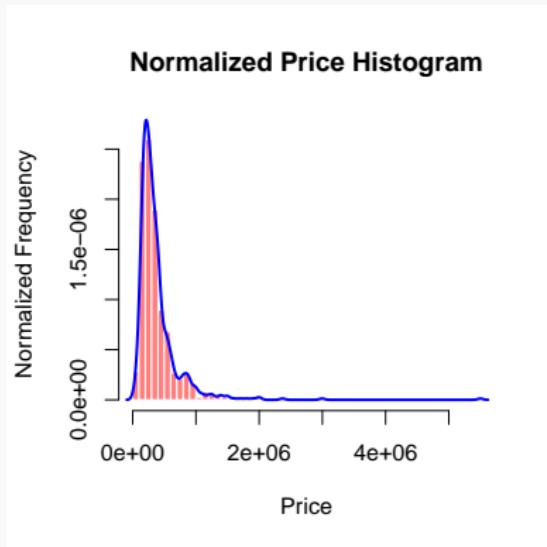
1. a **discrete RV** takes on a finite or countable number of values.

Ex: The number of bedrooms, B , of a home in our dataset is a random variable. B is discrete.

2. a **continuous RV** usually takes on all values in some range (a, b) .

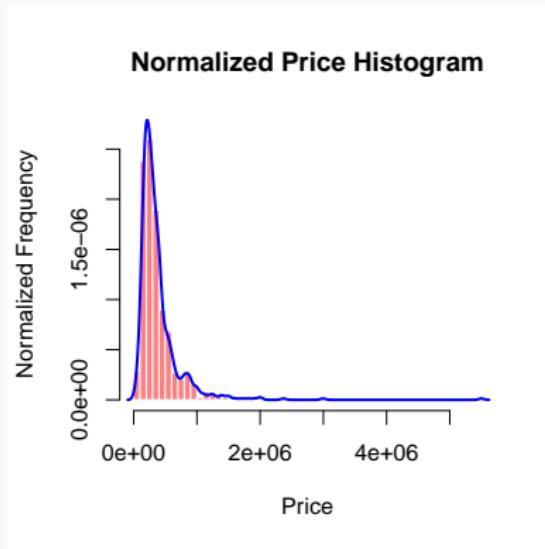
Ex: The observed price, P , of a home given the square footage is a random variable. P is continuous (can take on all values between 0 and ∞).

Probability Distributions



The **probability distribution** of a continuous RV X is given by a function, $p(X)$. The area under p over (a, b) describes the probability of observing values between $X = a$ and $X = b$. p is called the **probability density function** (pdf) of X .

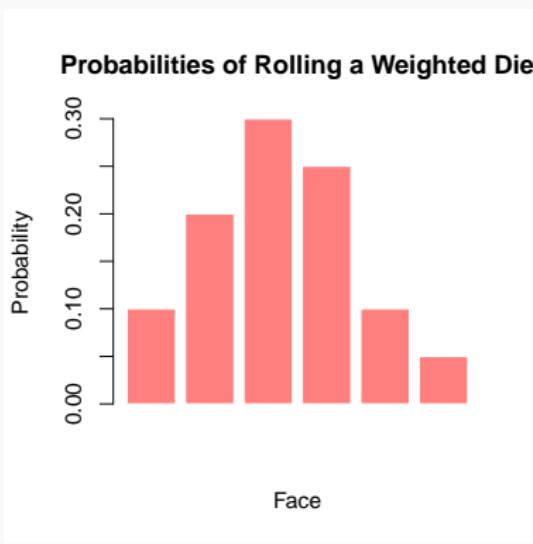
Probability Distributions



The pdf can provide intuition for how the RV behaves.

For example, the pdf gives us a sense of which values are more likely to be observed compared to others.

Probability Distributions



The **probability distribution** of a discrete RV X is given a function, $p(X)$. $p(a)$, written $p(X = a)$, is the probability of observing $X = a$.

p is called the **probability density function** (pdf) or the **probability mass function** (pmf) of X .

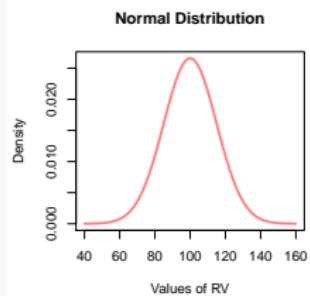
Descriptive Statistics

For many distributions, we can completely describe the shape of the pdf using just a few quantities. These quantities are usually:

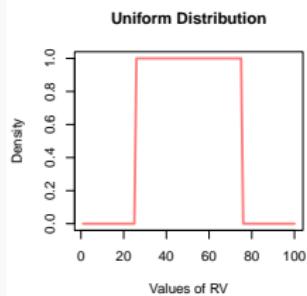
1. **(Measuring the ‘center’)** The **mean** measures the average of the outcomes, weighted by how likely is each outcome. The **median** divides the area under the pdf into two equal parts.
2. **(Measuring the ‘peak’)** The **mode** is the outcome that is the most likely (gives the highest value for the pdf).
3. **(Measuring the ‘spread’)** The **variance** is measures the average difference between outcomes and the mean, weighted by how likely is each outcome.

Common Types of Distributions (Continuous)

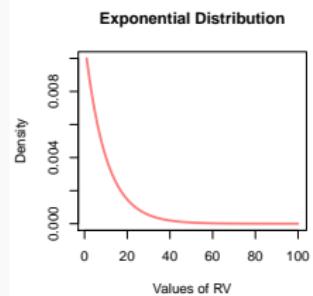
Normal Distribution



Uniform Distribution



Exponential Distribution



$$p(X) = \frac{1}{\sqrt{2\sigma^2}\pi} \exp\left\{-\frac{(X-\mu)^2}{2\sigma^2}\right\}$$

$$\mu \in \mathbb{R}, \sigma > 0$$

$$X \in (-\infty, \infty)$$

$$X \sim \mathcal{N}(\mu, \sigma)$$

$$p(X) = \begin{cases} \frac{1}{b-a}, & a \leq X \leq b \\ 0, & \text{otherwise} \end{cases}$$

$$a, b \in \mathbb{R}$$

$$X \in [a, b]$$

$$X \sim U(a, b)$$

$$p(X) = \lambda e^{-\lambda X}$$

$$\lambda > 0$$

$$X \in [0, \infty)$$

$$X \sim \text{Exp}(\lambda)$$

An Example

Recall that our model for observed housing prices, P , given square footage, A is

$$P = m \cdot A + b + \epsilon \quad (\text{model for observed prices})$$

where ϵ is a random noise variable.

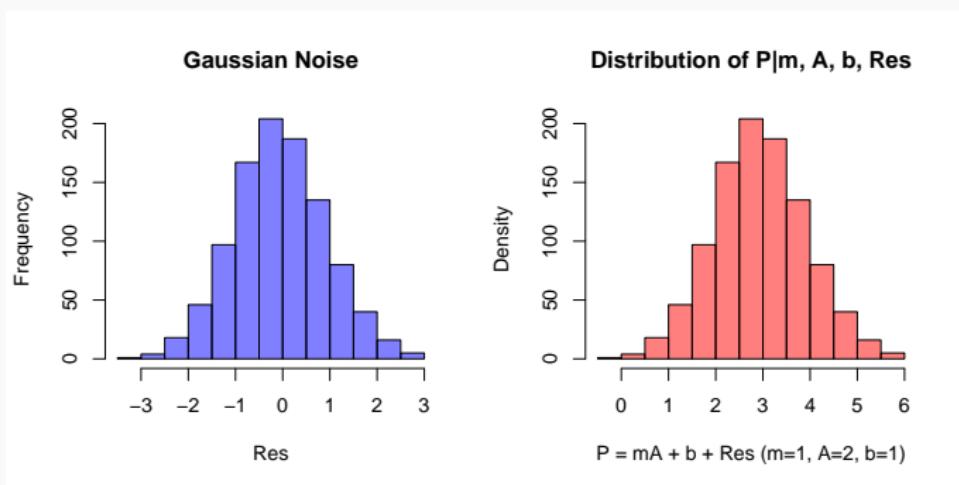
What kind of distributions $P|m, b, A, \epsilon$?

Try it yourself: Come up with a conjecture on the distribution of $P|m, b, A, \epsilon$ by generating values of ϵ and use them to generate values for P .

An Example

Let's fix $m = 1, b = 1, A = 2$ and choose a normal distribution for ϵ , say $\epsilon \sim \mathcal{N}(0, 1)$.

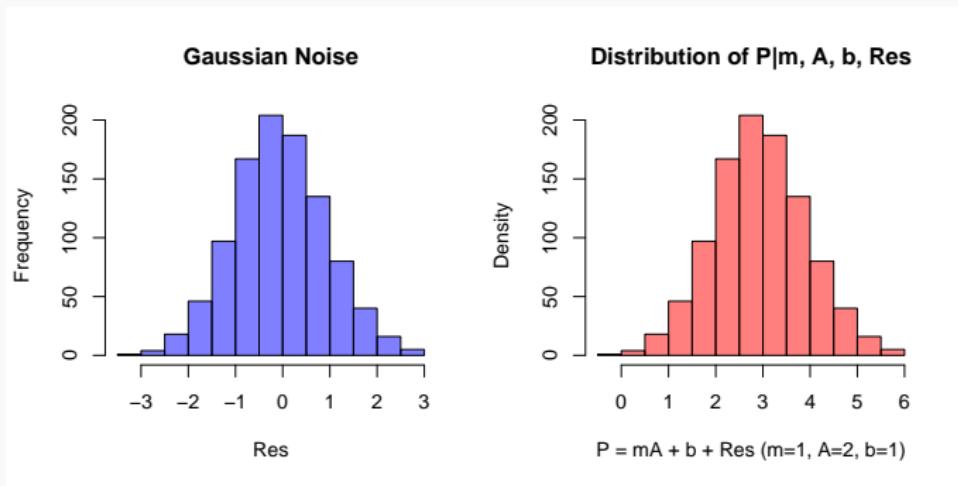
What is the distribution of $P|m, b, A, \epsilon?$



An Example

Let's fix $m = 1, b = 1, A = 2$ and choose a normal distribution for ϵ , say $\epsilon \sim \mathcal{N}(0, 1)$.

What is the distribution of $P|m, b, A, \epsilon$?



$P|m, b, A, \epsilon$ looks normally distributed:

$$P|m, b, A, \epsilon \sim \mathcal{N}(m \cdot A + b, 1)$$

Linear Regression (Univariate)

Back to Our Linear Model

Recall that our statistical model for observed housing prices is

$$P = m \cdot A + b + \epsilon$$

Suppose that $\epsilon \sim \mathcal{N}(0, \sigma^2)$, then

$$P|m, b, A, \epsilon \sim \mathcal{N}(m \cdot A + b, \sigma^2).$$

Can we use this information to find the best linear model for the observed data?

Try it yourself: Use python to calculate $p(P|m, b, A, \epsilon)$ for a couple of choices of m and b .

Back to Our Linear Model

Let's find m_{MLE} and b_{MLE} so that $p(P|m_{MLE}, b_{MLE}, A, \epsilon)$ is maximal.

That is, the model $P = m_{MLE} \cdot A + b_{MLE} + \epsilon$ explains the observed data with the highest probability.

The above model is called the **maximum likelihood estimator (MLE)**.

Ordinary Least Squares

For Gaussian noise, $\epsilon \sim \mathcal{N}(0, 1)$, finding the MLE model is the same as minimizing the MSE loss function!

Recall, given a set of points $(A_1, P_1), \dots, (A_n, P_n)$, the MSE of a linear model $P = mA + b$ is

$$L(m, b) = \sum_{i=1}^n |P_i - (mA_i + b)|^2$$

L is the lost (or cost) function. Our goal is to **find** m_{MLE} and b_{MLE} such that the lost, $L(m_{MLE}, b_{MLE})$, is **minimal**:

$$(m_{MLE}, b_{MLE}) = \operatorname{argmin}_{m,b} L(m, b).$$

Evaluating Models

Given a set of home listings, we can now fit a maximum likelihood linear model

$$P = m_{MLE} \cdot A + b_{MLE}$$

by minimizing the sum of squared residuals (ordinary least squares or OLS). The ‘error’ made by our model in fitting the data is

$$L(m_{MLE}, b_{MLE}) = \sum_{i=1}^n |P_i - (m_{MLE} \cdot A_i + b_{MLE})|^2.$$

This is called the **training error**.

But we also need to evaluate our model on new data that it has not yet seen, **test data**.

Evaluating Models

The general process of model fitting and evaluation:

1. Given a set of data $(A_1, P_1), \dots, (A_n, P_n)$, split the data into a **training set** and a **test set**.
2. Fit the model on the training set, report error
3. Fit the model on the testing set, report error

Linear Regression (Multivariate)

Linear Regression in Multiple Variables

It's a bit unreasonable for price of a home to depend on square footage alone. In reality, P most likely depends on some combination of square footage, A , number of bedrooms Bd and the number of bathrooms Ba .

The easiest relationship between all 4 variables is again linear

$$P = a_0 + a_1A + a_2Bd + a_3Ba + \epsilon.$$

Again, if we take $\epsilon \sim \mathcal{N}(0, 1)$, P (given the model) is a random variable with a normal distribution,
 $P \sim \mathcal{N}(a_0 + a_1A + a_2Bd + a_3Ba, 1)$.

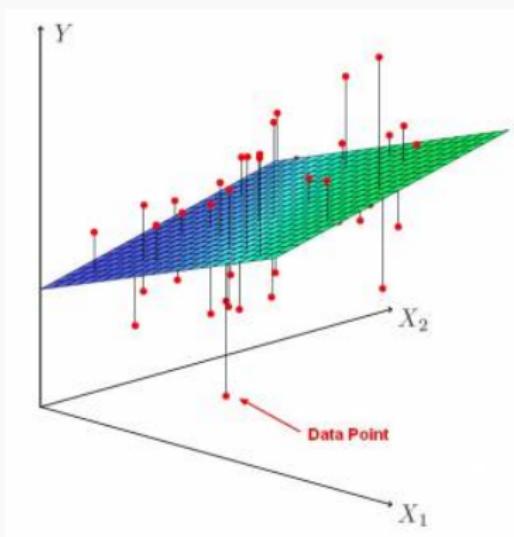
Just like before, the values of a_0, \dots, a_3 which maximizes the likelihood of the data (MLE model) can be found by minimizing sum of squared residuals.

Linear Regression in Multiple Variables

When we are fitting a multi-linear regression model

$$P = a_0 + a_1 A + a_2 Bd + a_3 Ba + \epsilon.$$

to our data. Instead of finding the best fitting line, we find the best fitting **plane**:



Try It Yourself!

Try to fit multi-linear models with a few sets of variables other than square footage and report the training and testing errors.

Polynomial Regression

Polynomial Regression

As we've noticed, our linear models (univariate and multivariate) don't seem to fit the housing data very well.

Maybe this is because the underlying relationship between price and square footage (or number of rooms) isn't linear. Perhaps the model we want is polynomial

$$P = a_0 + a_1 A + a_2 A^2 + \epsilon, \quad \epsilon \sim \mathcal{N}(0, 1)$$

Note: our polynomial model is secretly a multi-linear model with two variables, A and A^2 .

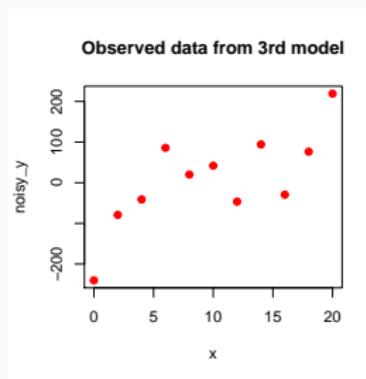
Try It Yourself!

Try to fit polynomial model with a few different degrees and report the training and testing errors.

Overfitting

One thing we notice immediately was that the error on the training set decreases as the degree of the polynomial increases.

So maybe we generally want to pick very high degree polynomials to model our data?

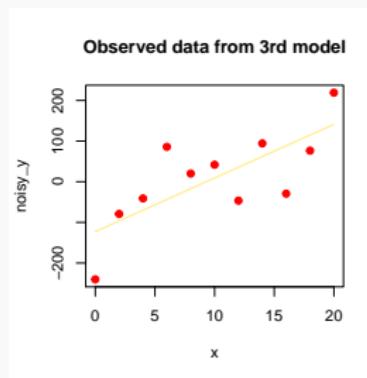


What is happening to our model as the degree increases?

Overfitting

One thing we notice immediately was that the error on the training set decreases as the degree of the polynomial increases.

So maybe we generally want to pick very high degree polynomials to model our data?

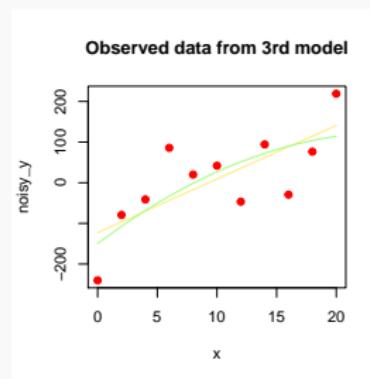


What is happening to our model as the degree increases?

Overfitting

One thing we notice immediately was that the error on the training set decreases as the degree of the polynomial increases.

So maybe we generally want to pick very high degree polynomials to model our data?

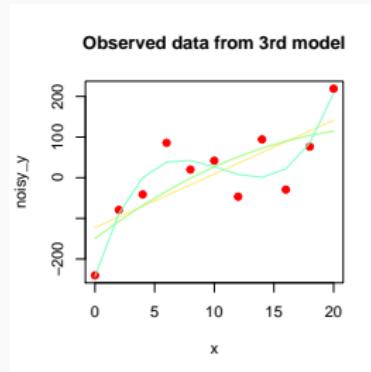


What is happening to our model as the degree increases?

Overfitting

One thing we notice immediately was that the error on the training set decreases as the degree of the polynomial increases.

So maybe we generally want to pick very high degree polynomials to model our data?

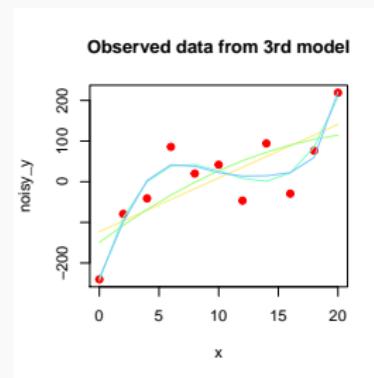


What is happening to our model as the degree increases?

Overfitting

One thing we notice immediately was that the error on the training set decreases as the degree of the polynomial increases.

So maybe we generally want to pick very high degree polynomials to model our data?

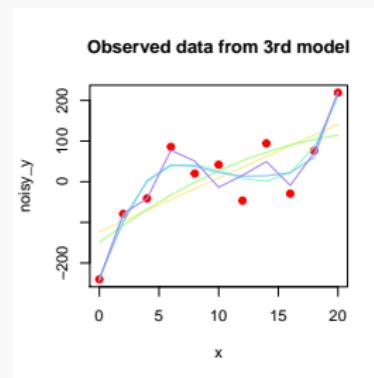


What is happening to our model as the degree increases?

Overfitting

One thing we notice immediately was that the error on the training set decreases as the degree of the polynomial increases.

So maybe we generally want to pick very high degree polynomials to model our data?

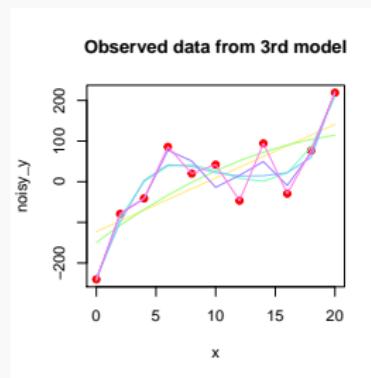


What is happening to our model as the degree increases?

Overfitting

One thing we notice immediately was that the error on the training set decreases as the degree of the polynomial increases.

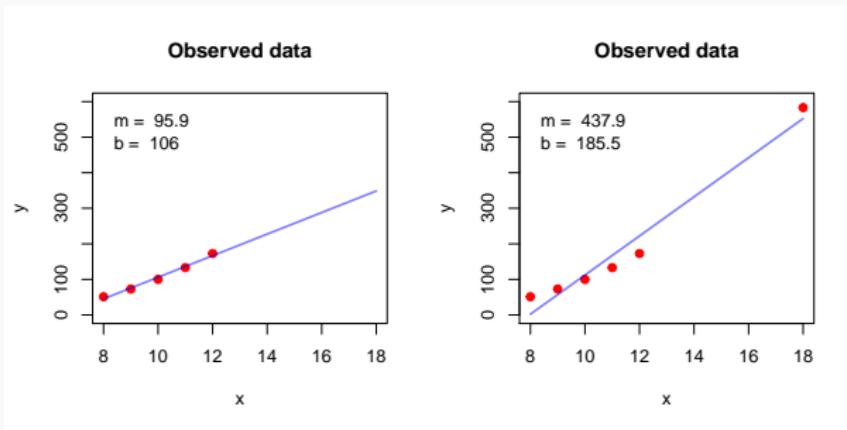
So maybe we generally want to pick very high degree polynomials to model our data?



What is happening to our model as the degree increases?

Overfitting

Overfitting can happen with linear regression too!



In multiple linear regression, what happens when we have N number of observations and N number of explanatory variables?

Overfitting

Overfitting happens when we learn parameters or rules that are too specific to the training set, so much that our model is not useful in explaining new data (we do great on train data but poorly on test).

Overfitting can happen when we have too few observations compared to the number of variables in our model with which we try to explain the observations.

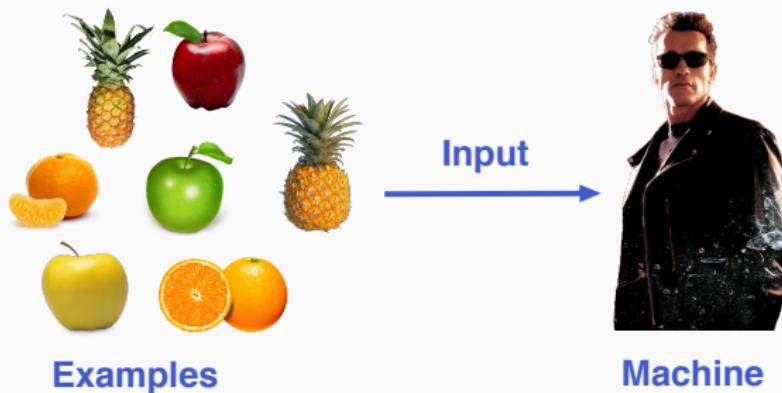
Later, we'll see that overfitting can be curbed by regularization and variable selection.

What is Machine Learning?

Intuition

The goal of machine learning is to be able to teach machines to make decisions based on previous experience, just like humans.

For example, we give the machine examples of a type of object or scenario.



Intuition

The goal of machine learning is to be able to teach machines to make decisions based on previous experience, just like humans.

For example, we give the machine examples of a type of object or scenario. We hope that machine will recognize a new instance of that type of object or scenario.

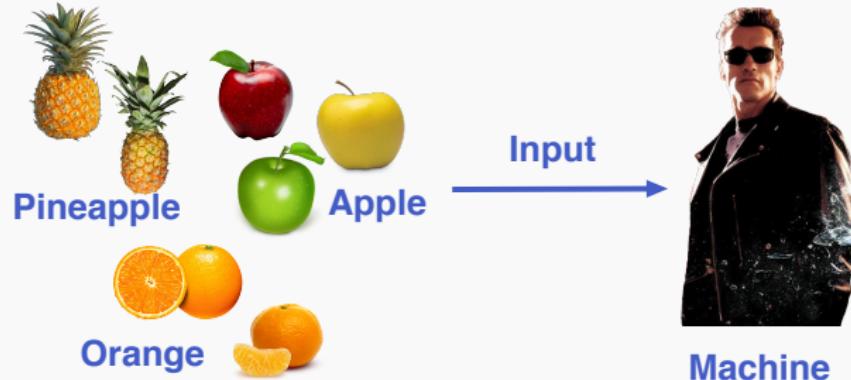


Intuition

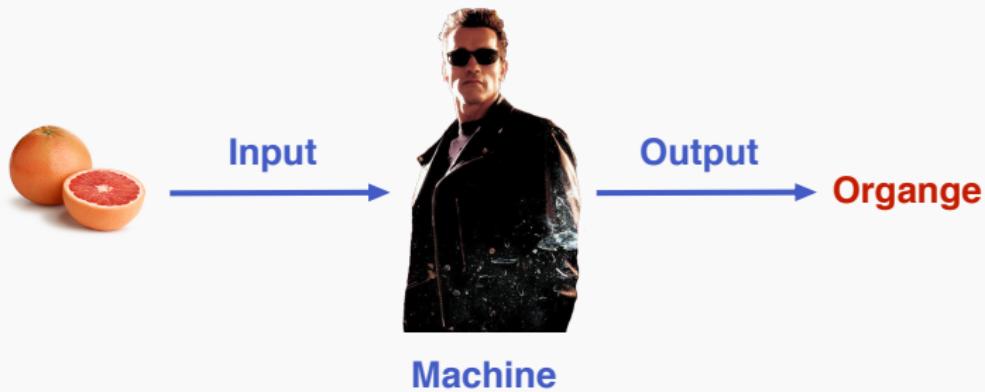
The goal of machine learning is to be able to teach machines to make decisions based on previous experience, just like humans.

The ‘learning’ done by the machine is usually fitting a statistical model to a set of training data. The machine can then use this calibrated model to make decisions when encountering new data.

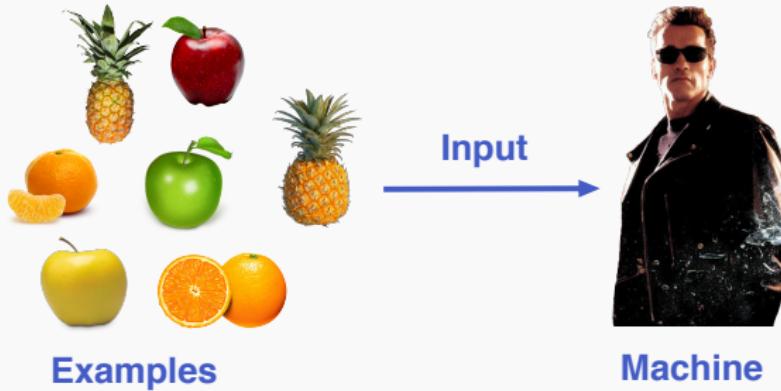
Supervised Learning



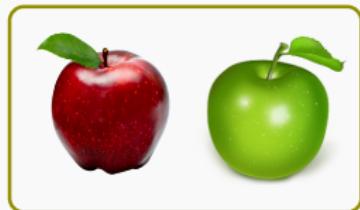
Supervised Learning



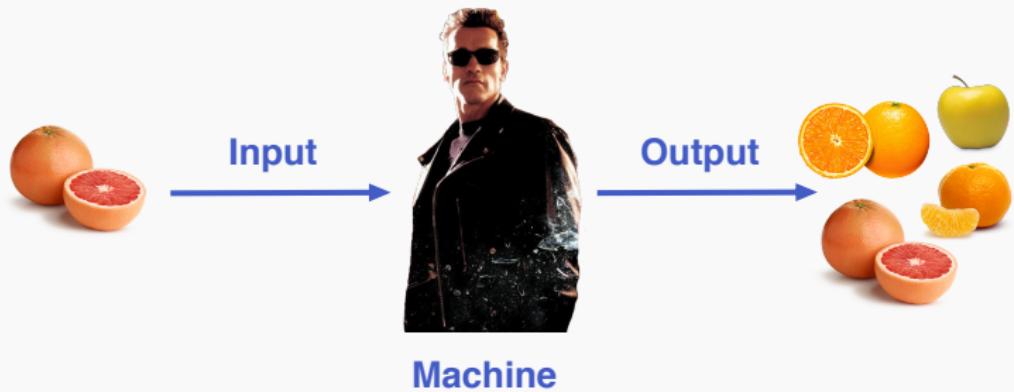
Unsupervised Learning



Unsupervised Learning



Unsupervised Learning



Final Thoughts

Question: What is the difference between machine learning and data science?

Question: Is regression (linear or polynomial) supervised or unsupervised learning? Why?

Question: What might be the pros and cons of supervised (resp unsupervised learning)? Why?

Question: Why is it important to use statistical models in representing data?