

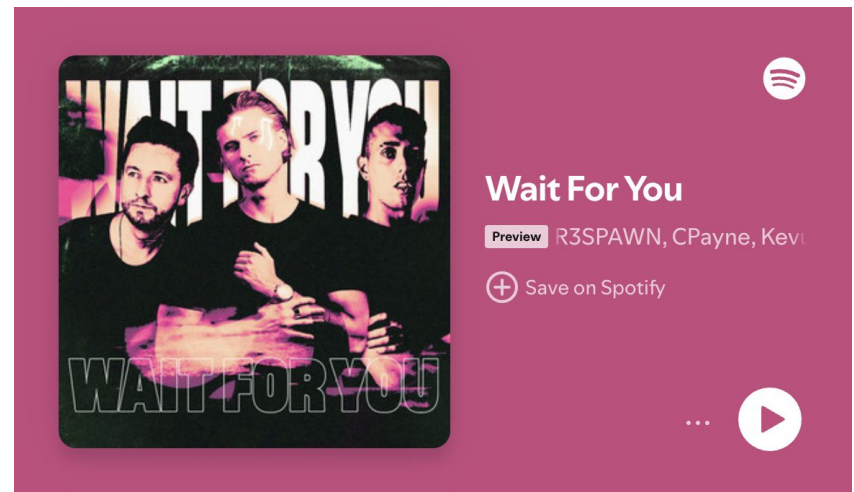
PROFESSIONAL CERTIFICATE IN MACHINE LEARNING AND ARTIFICIAL INTELLIGENCE

Let's give everyone a couple of minutes to join...

Module 8

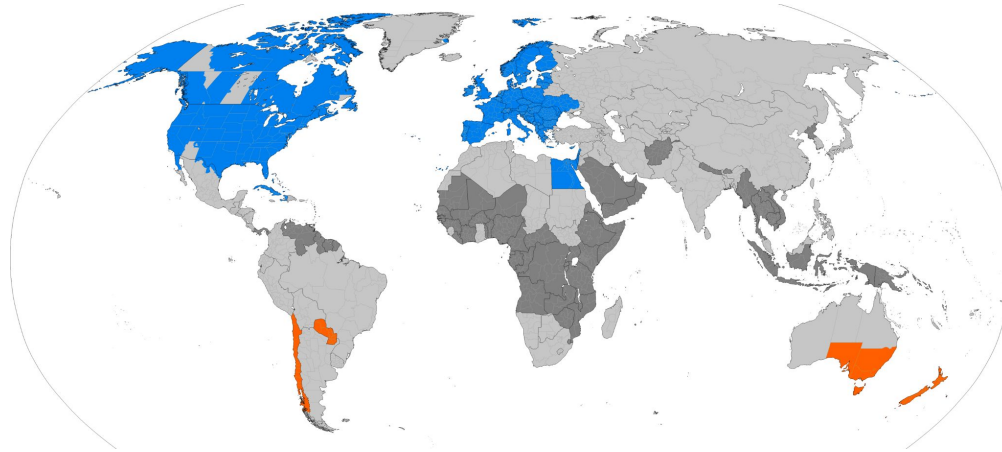
Feature Engineering and Overfitting

Office Hours with Viviana Márquez
October 24, 2024



<https://open.spotify.com/track/63X08jZU1piJix8kvfoaq2?si=8e87e15cba5045b2>

Daylight saving time!



Always check Canvas for the most up-to-date information regarding office hours!

Tool to convert to your timezone: <https://www.worldtimebuddy.com/>



Europe
DST ends



North America
DST ends

AGENDA

- Required activities for Module 8
- Content review Module 8: Feature Engineering and Overfitting
- Code
- Questions


AGENDA

- Required activities for Module 8
- Content review Module 8: Feature Engineering and Overfitting
- Code
- Questions

Required Activities for Module 8

- Knowledge Check 8.1: Parabolic Model Fitting and Nonlinear Features
- Knowledge Check 8.2: Scikit-Learn Transformers
- Knowledge Check 8.3: Scikit-Learn Pipelines
- Codio Assignment 8.1: Scikit-Learn Pipeline
- Check 8.4: Order 0 Through 6 Models on Vehicle Data
- Knowledge Check 8.5: The Dangers of Overfitting
- Codio Assignment 8.2: Comparing Complexity and Variance
- Check 8.6: Overfitting and Validation
- Knowledge Check 8.7: Test Sets
- Codio Assignment 8.3: Evaluating Multiple Models

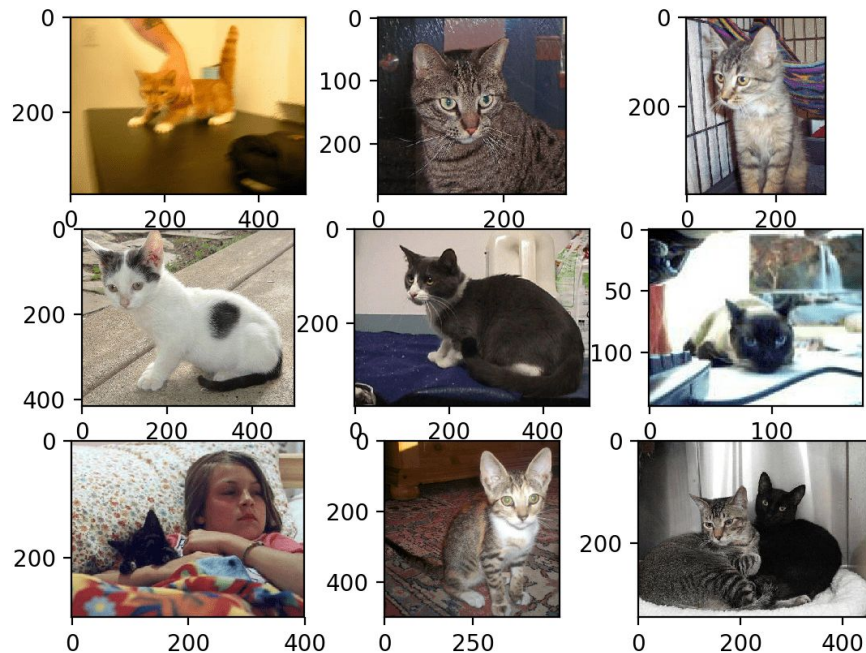
AGENDA

-  Required activities for Module 8
- Content review Module 8: Feature Engineering and Overfitting
- Code
- Questions

Content review Module 8: Feature Engineering and Overfitting

- Quick recap Linear Regression
- Train/val/test datasets
- Sources of error in a model
- Bias-variance tradeoff
- Feature Engineering

What is Machine Learning?



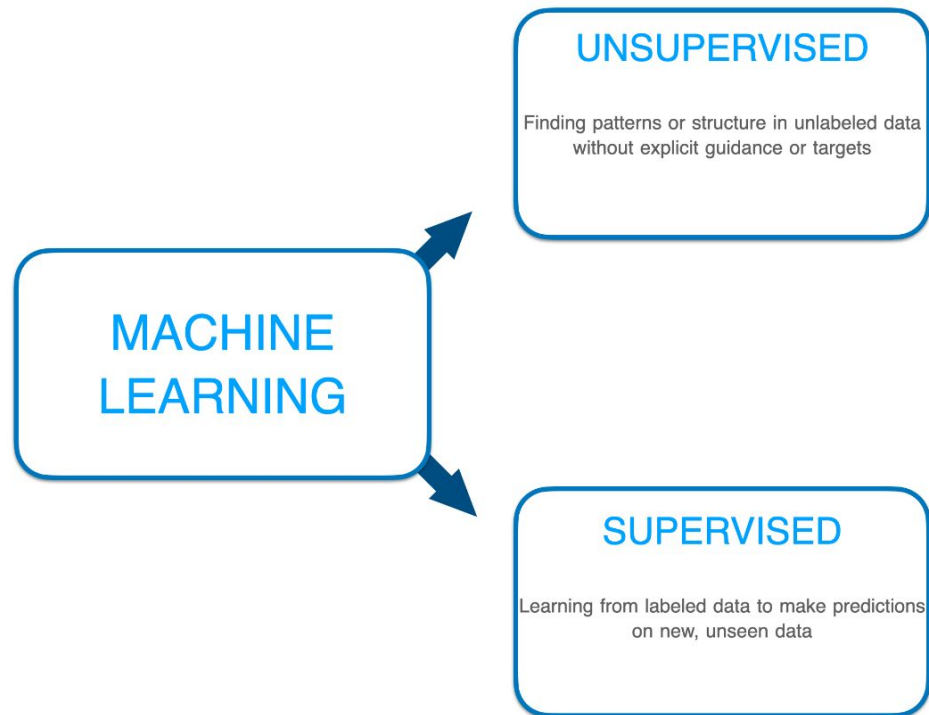
In machine learning, instead of explicitly programming a computer with specific instructions to perform a task, we provide it with large amounts of data and allow it to learn how to **generalize** from that data.

The Machine Learning landscape

Which tool would you use to hit a nail?



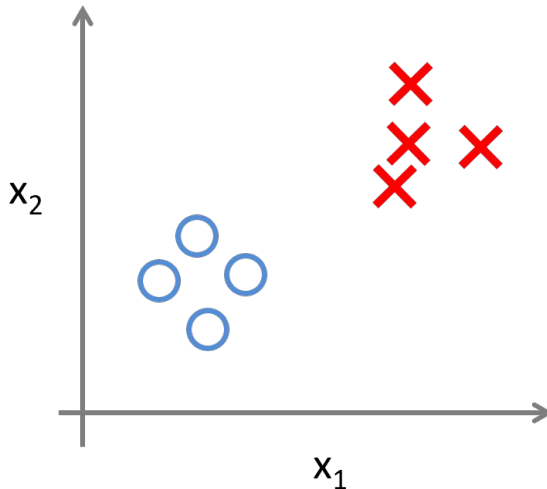
First Question: Do we have labels?



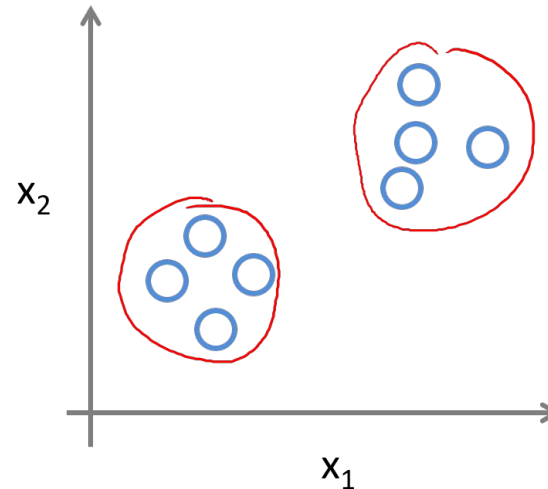
The Machine Learning landscape

Supervised learning vs Unsupervised learning

Supervised Learning



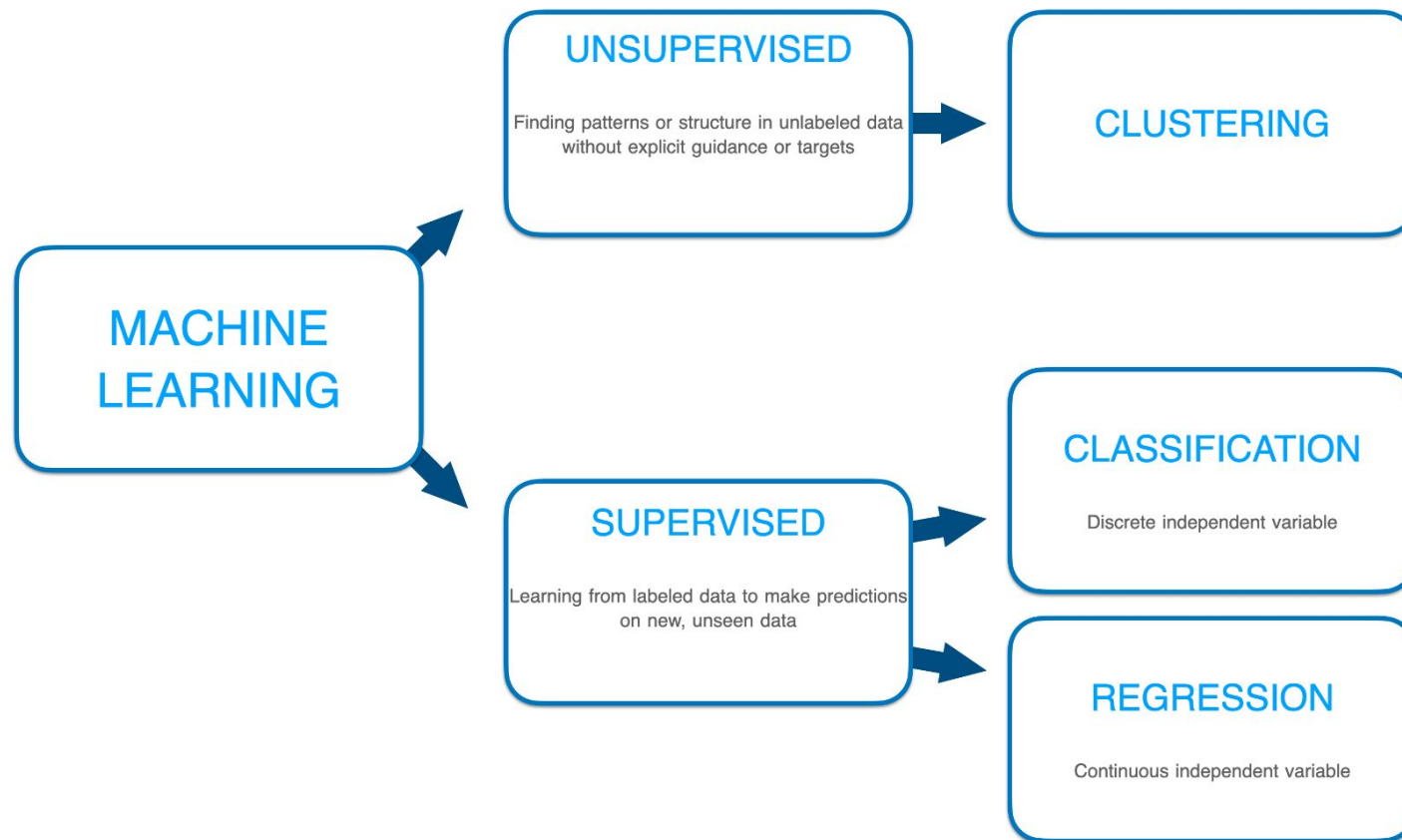
Unsupervised Learning



- **Supervised learning:** Problems with labels. Its aim is to predict data based on the labeled information.
- **Unsupervised learning:** Problems without labels. Its goal is to uncover patterns, structures, and relationships.

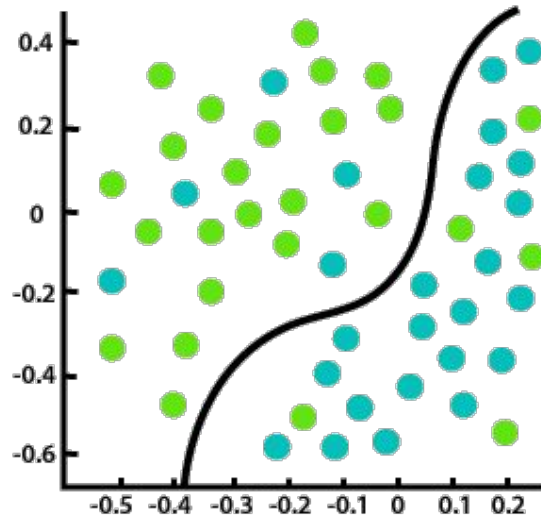
First Question: Do we have labels?

Second Question: Are our labels categorical or numerical?

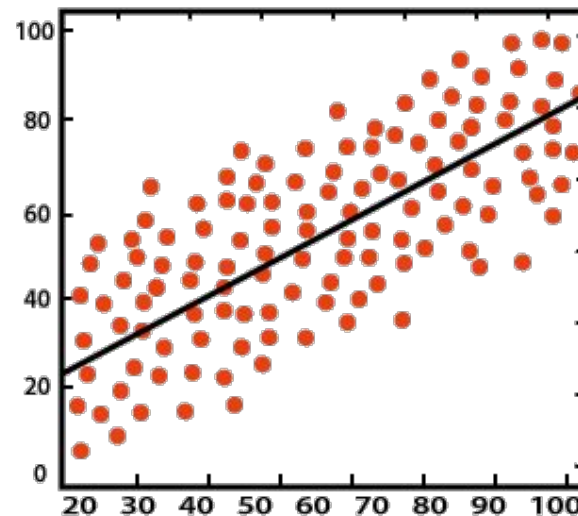


The Machine Learning landscape

Regression vs Classification



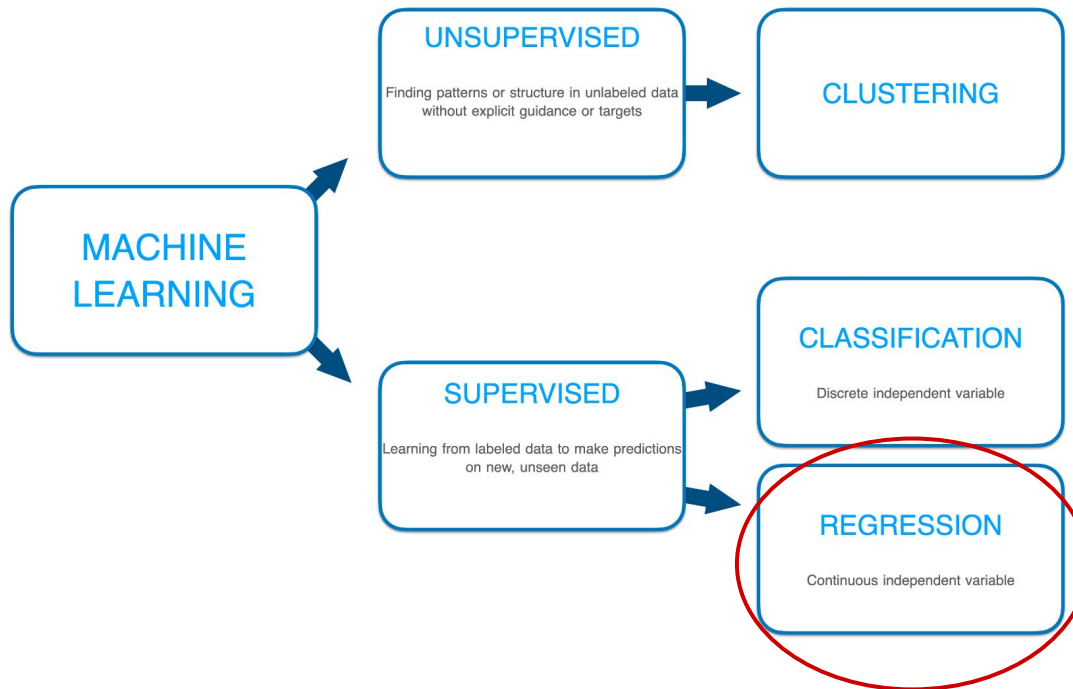
Classification



Regression

- **Regression:** Quantitative (continuous/numerical) target variable
- **Classification:** Qualitative (discrete/categorical) target variable

What type of data do we have?

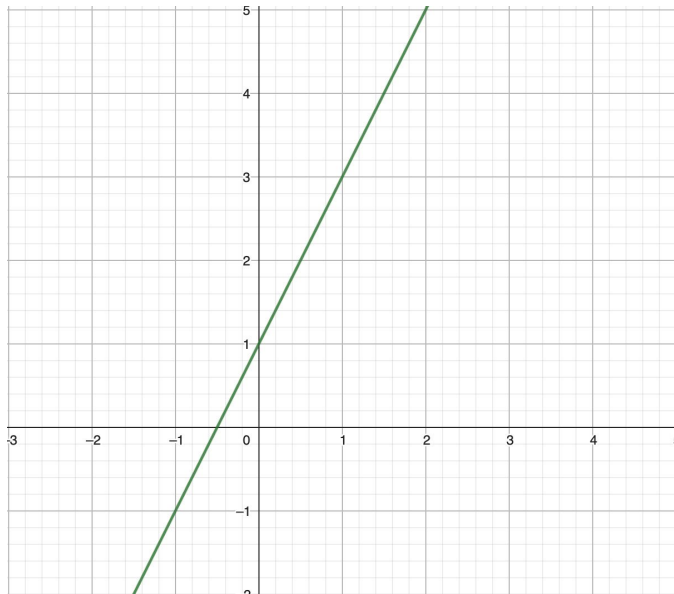


Linear Regression

- Linear regression is one of the most well known machine learning models
- It is:
 - **Supervised** (has labels)
 - **Regression** (the labels are numerical values)
- The goal of a linear regression is to model the relationship between a dependent variable (target) and one or more independent variables (features) by fitting a linear equation

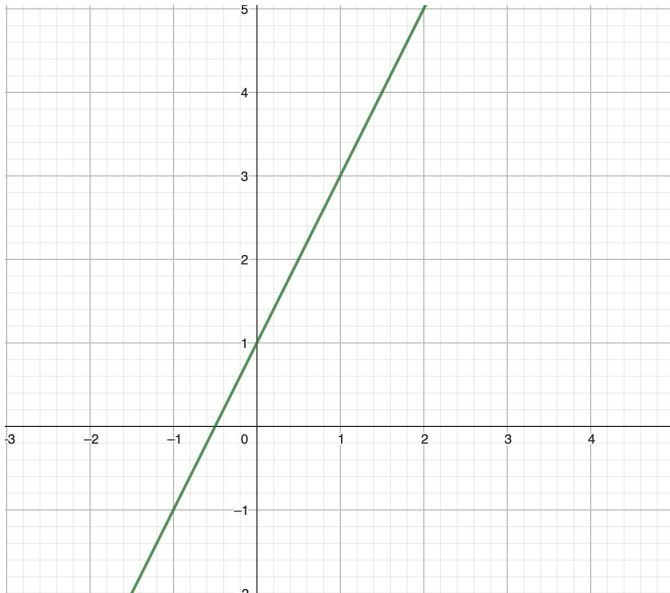
A tiny refresher

- What is the general equation of a straight line?



A tiny refresher

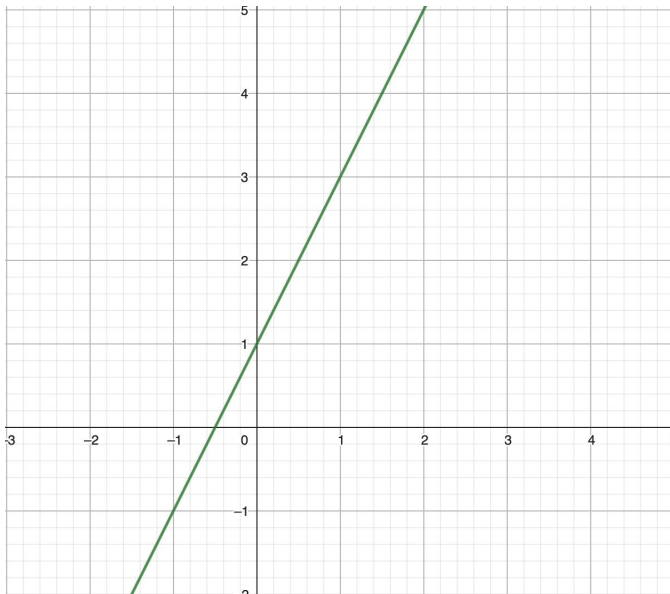
- What is the general equation of a straight line?



$$y = mx + b$$

A tiny refresher

- What is the general equation of a straight line?

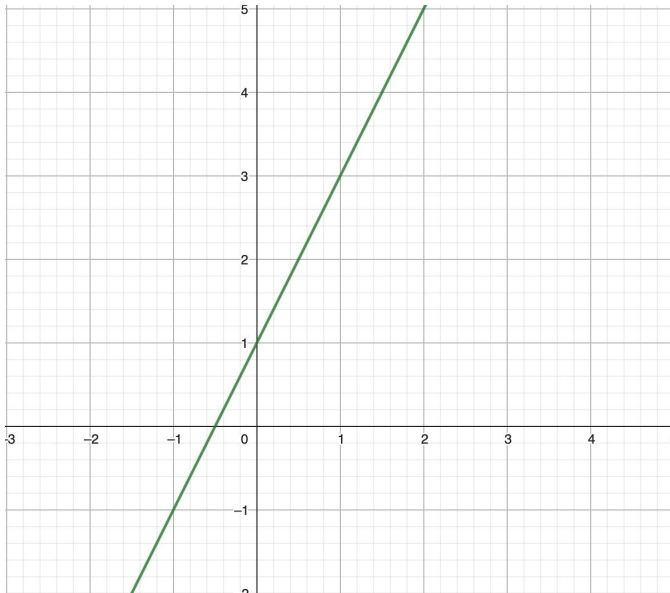


$$y = mx + b$$

Slope
(how steep the line is)

A tiny refresher

- What is the general equation of a straight line?



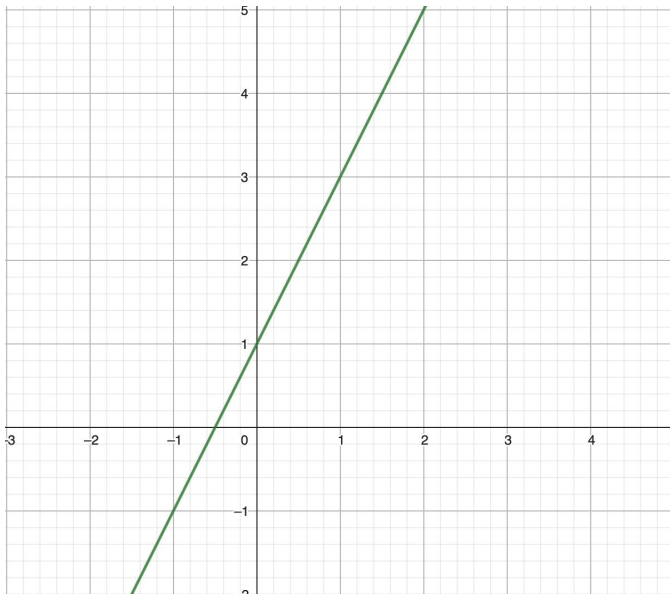
$$y = mx + b$$

Slope
(how steep the line is)

y-intercept

A tiny refresher

- What is the general equation of a straight line?



$$y = mx + b$$

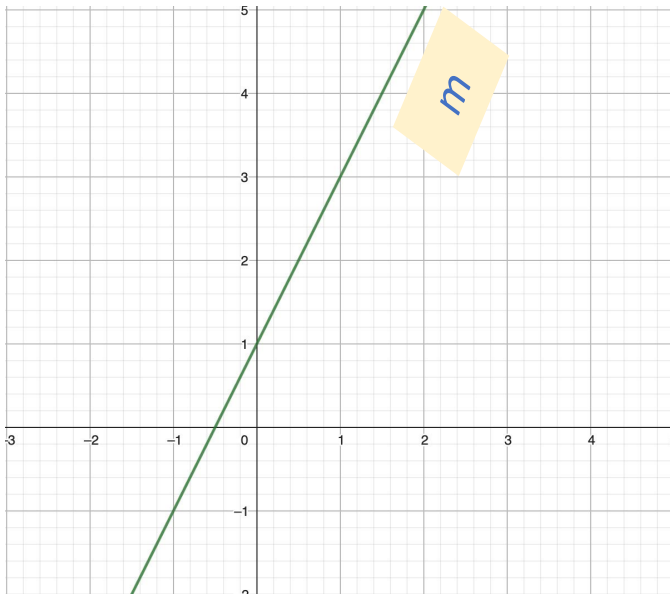
Slope
(how steep the line is)

y-intercept

$$y =$$

A tiny refresher

- What is the general equation of a straight line?



$$y = mx + b$$

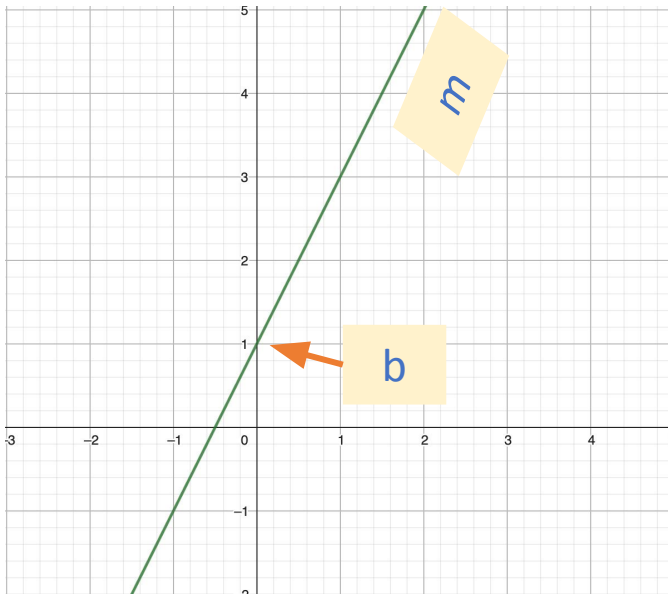
Slope
(how steep the line is)

y-intercept

$$y = 2x$$

A tiny refresher

- What is the general equation of a straight line?



$$y = mx + b$$

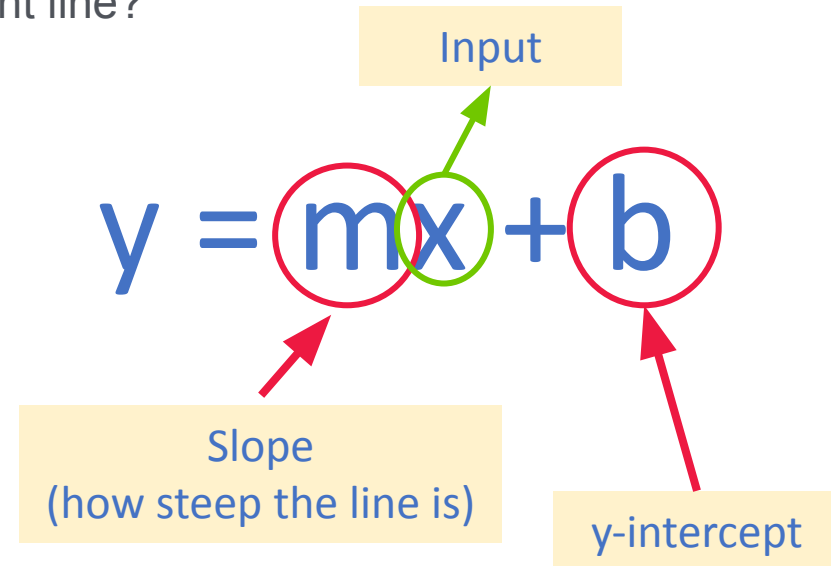
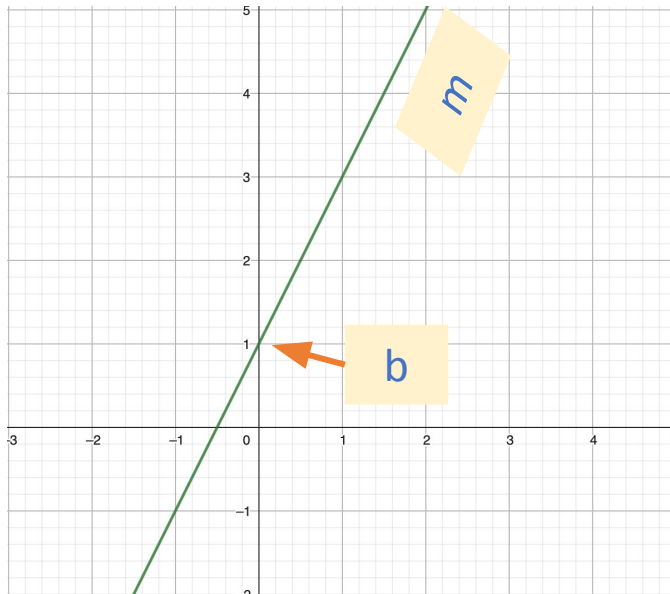
Slope
(how steep the line is)

y-intercept

$$y = 2x + 1$$

A tiny refresher

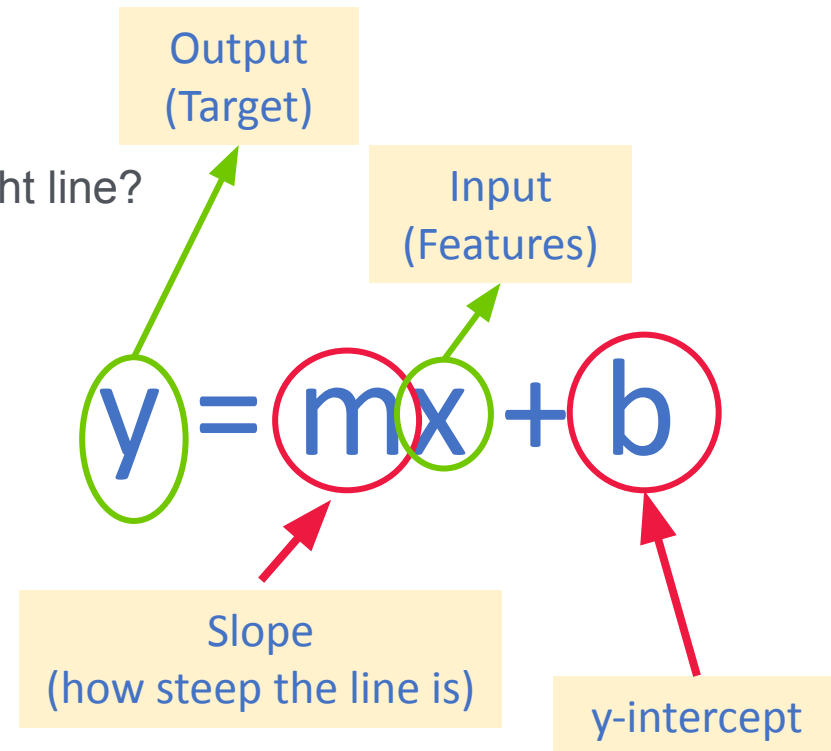
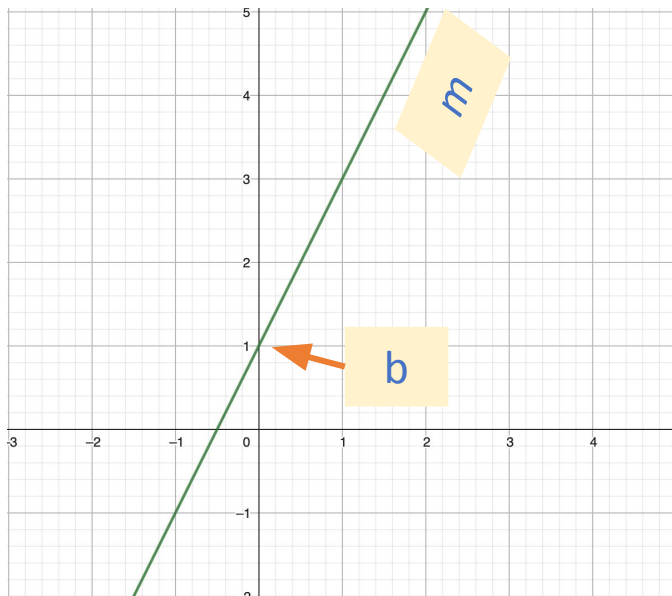
- What is the general equation of a straight line?



$$y = 2x + 1$$

A tiny refresher

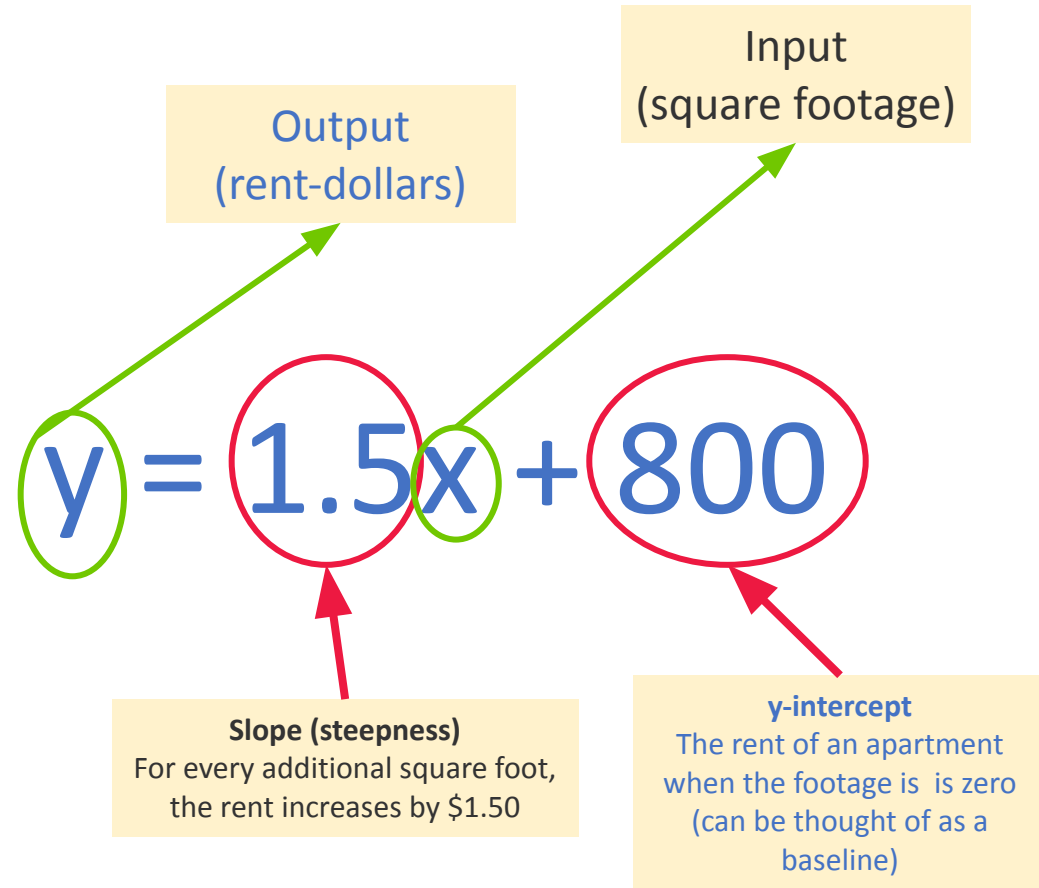
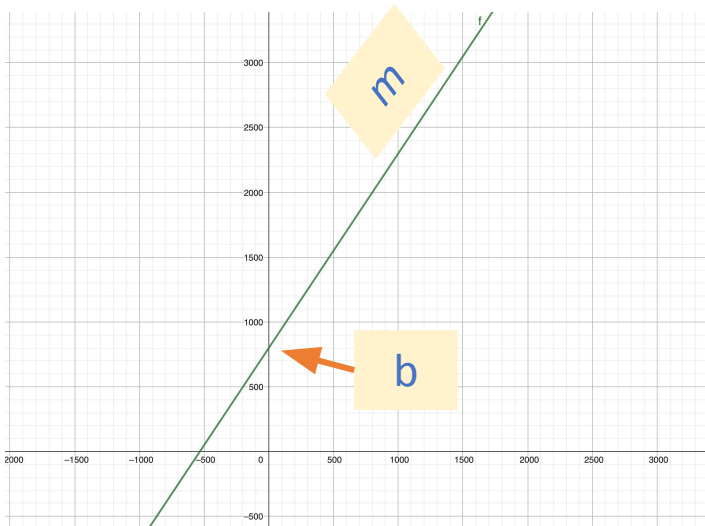
- What is the general equation of a straight line?



$$y = 2x + 1$$

A tiny refresher

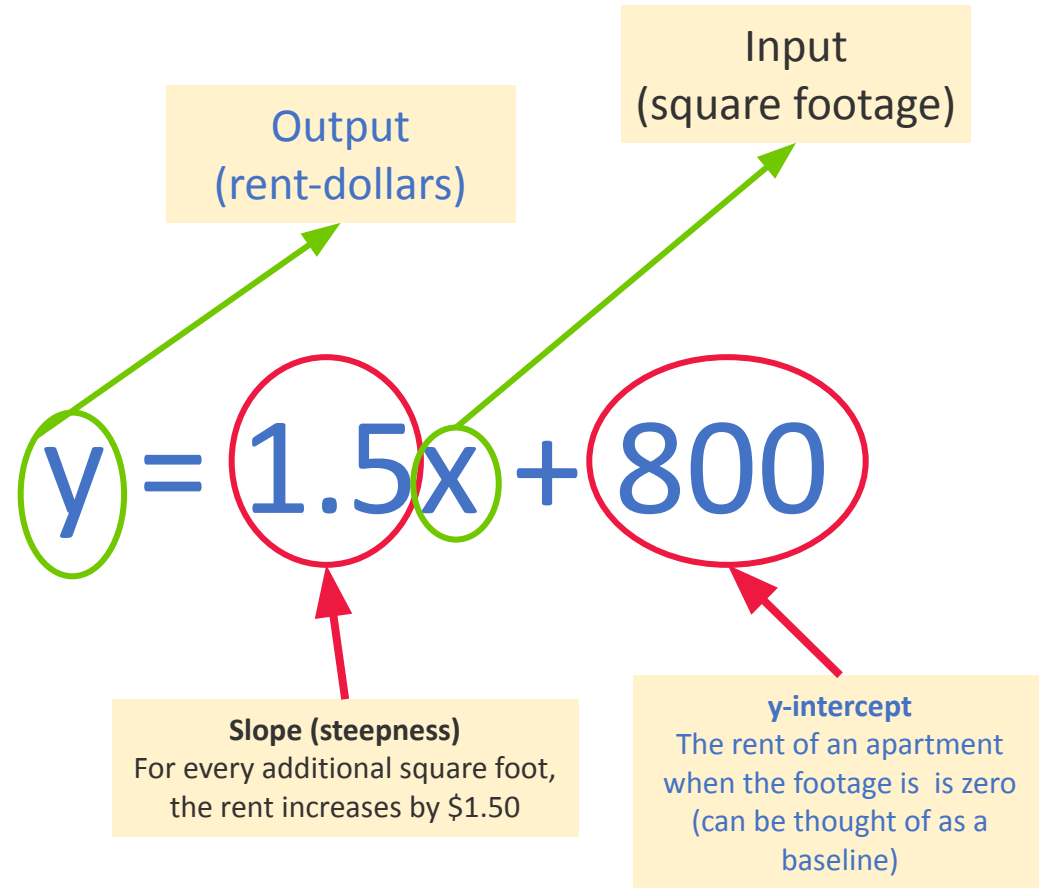
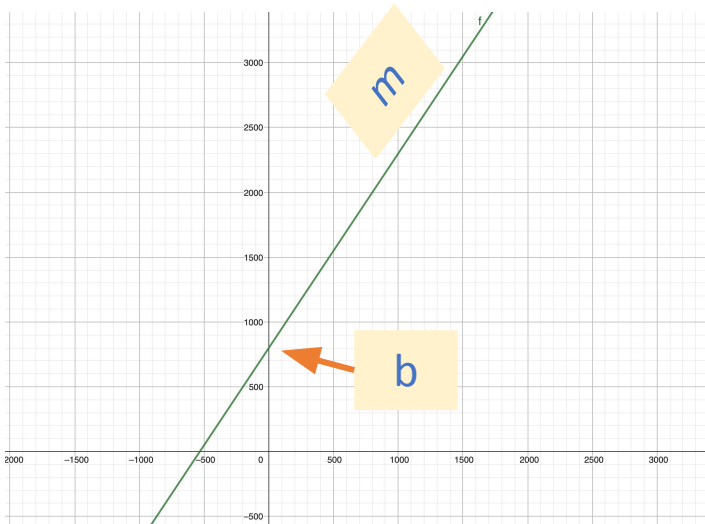
- Example: Predict monthly rent based on the sq ft of an apartment



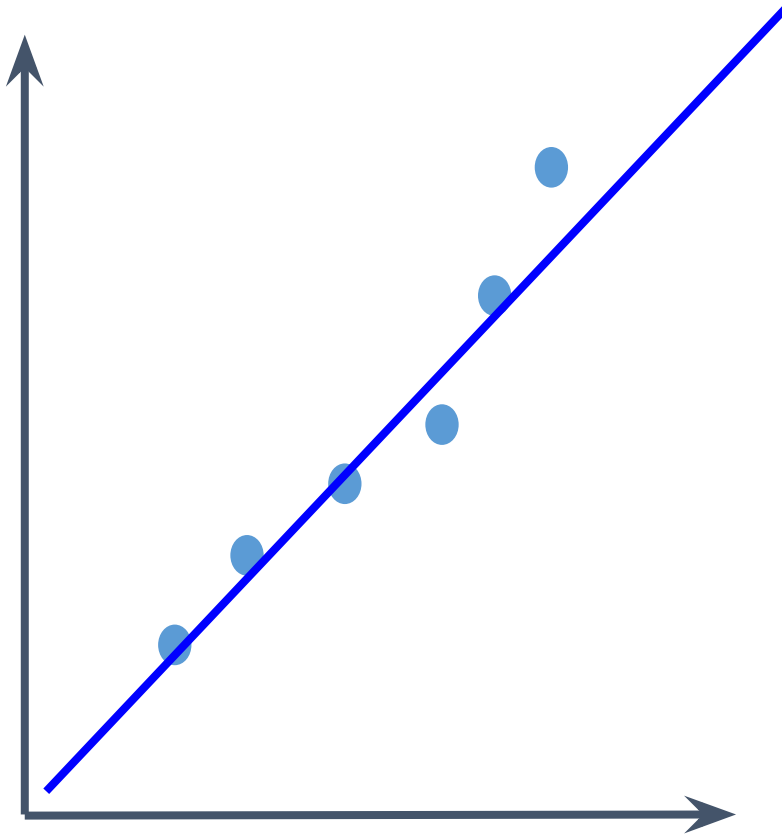
A tiny refresher

How can we find the best values for the slope and the y-intercept? 🤔

- Example: Predict monthly rent based on the sq ft of an apartment



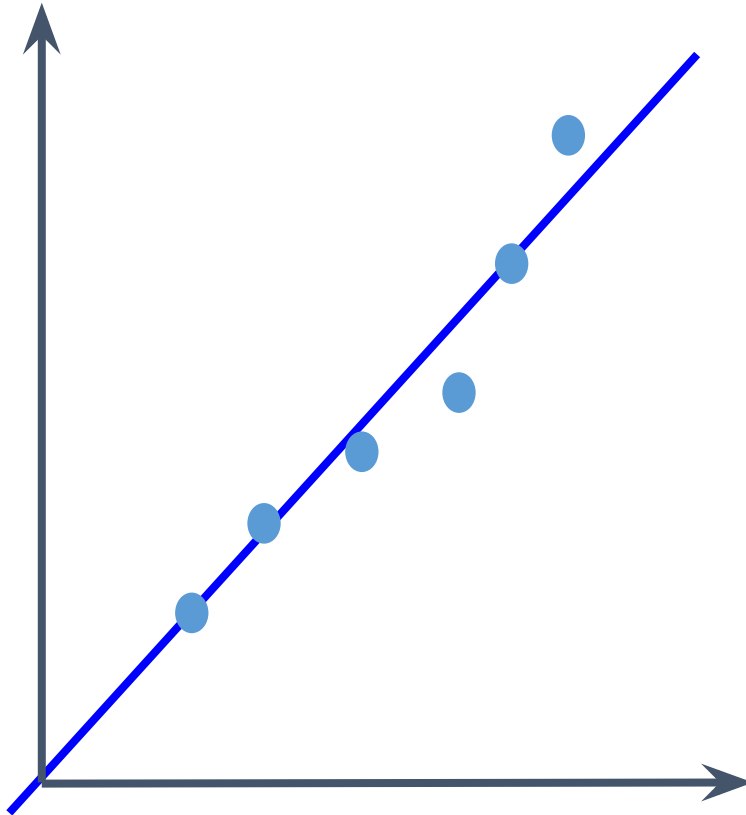
How can we find the best values for the slope and the y-intercept? 🤔



And that's Linear Regression!

- Linear regression is one of the most well known machine learning models
- It is:
 - **Supervised** (has labels)
 - **Regression** (the labels are numerical values)
- The algorithm tries to find the best-fitting straight line (in simple regression) or hyperplane (in multiple regression) that models the relationship between a dependent variable (target) and one or more independent variables (features)
- This "best fit" is often determined by minimizing the difference (or error) between the predicted values and the actual observed values, a process known as Ordinary Least Squares (OLS)
- Simple, interpretable, very fast, and the best for linear relationships
- Usually a lower bound on performance, but often form the foundation for other, more powerful techniques
- If we combine multiple linear models and add a twist to it (an activation function), we get a neural network; those are incredibly useful and powerful! 🔥🔥🔥

How can we find the best values for the slope and the y-intercept? 🤔



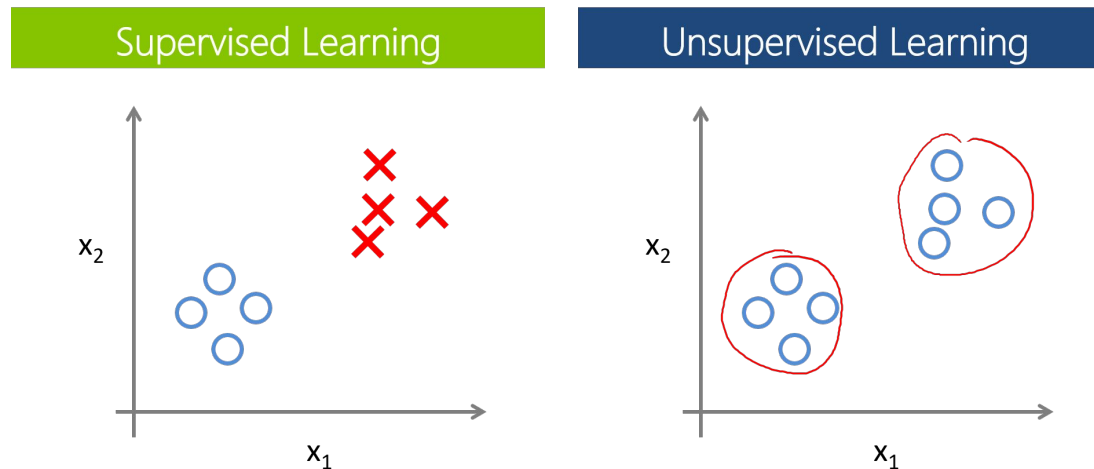
```
lm = LinearRegression()  
lm.fit(X, y)
```

It's so easy using code!



The Machine Learning landscape

Supervised learning vs Unsupervised learning



- **Supervised learning:** Problems with labels
- **Unsupervised learning:** Problems without labels



Train/val/test data sets



Is my model any good?

Train/val/test data sets

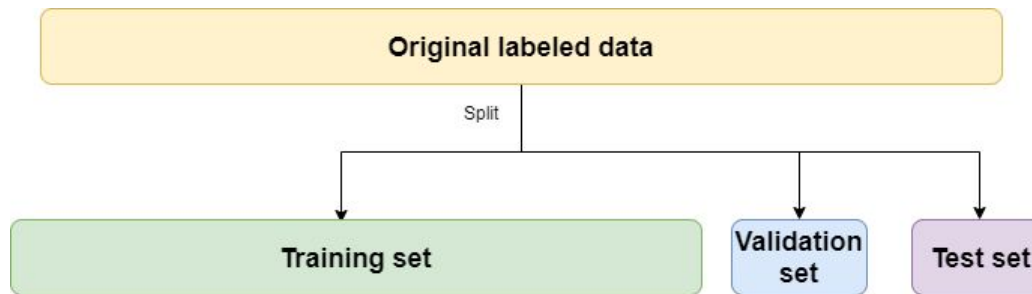


Is my model any good?

- A good model makes useful predictions on unknown, future data (it **generalizes**)
- It might not be *very* accurate, but if it's better than a coin flip, it might still be useful

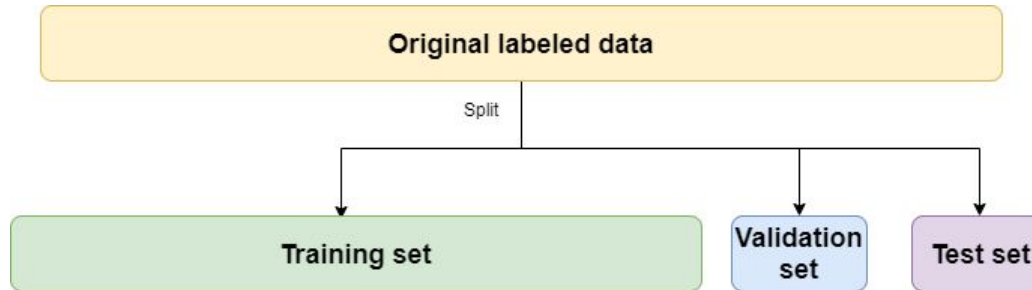
Train/val/test data sets

In order to make sure the model is learning to generalize, we split our dataset



- **Training data set:** Used to train machine learning model
- **Validation data set:** Used to tune the model hyperparameters
- **Test data set:** Evaluate the model's performance on unseen data, our best proxy on how the model will perform in the real world

Train/val/test data sets



- In general, putting 80% of the data in the training set, 10% in the validation set, and 10% in the test set is a good split to start with
- Most of the time, the data should be split randomly to be representative of the whole population

Train/val/test data sets



- It is important to split the data into train/val/test **BEFORE** doing any feature engineering or modeling to prevent **data leakage** and ensure that the model is truly evaluated on unseen data

Train/val/test data sets



- It is important to split the data into train/val/test **BEFORE** doing any feature engineering or modeling to prevent **data leakage** and ensure that the model is truly evaluated on unseen data

Data leakage: When information from val or test is used to inform the model during training or feature engineering

Train/val/test data sets



- It is important to split the data into train/val/test **BEFORE** doing any feature engineering or modeling to prevent **data leakage** and ensure that the model is truly evaluated on unseen data

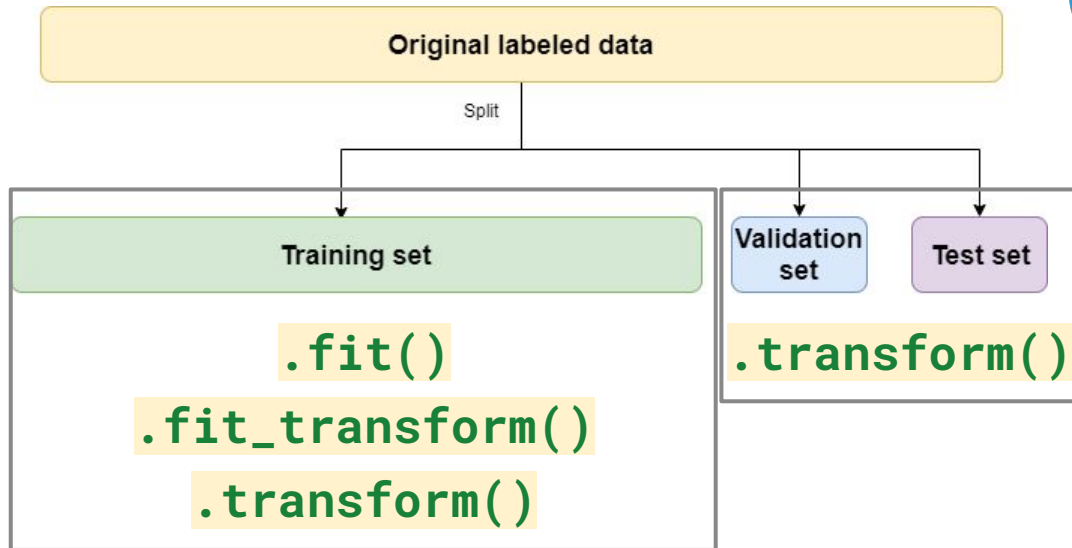
Data leakage: When information from val or test is used to inform the model during training or feature engineering



No peeking!!!

The test data set is used **ONLY** after you think you have the best model. It's the only true measure of generality.

Train/val/test data sets



- **.fit():**

This is the learning step, where the model or transformer observes the data to derive key parameters (e.g., coefficients in linear regression or decision boundaries in classification or the mean and variance in scaling). **ONLY** using on the train dataset. This is because the model should not learn from the val/test data—its purpose is to evaluate how well the model generalizes to unseen data.

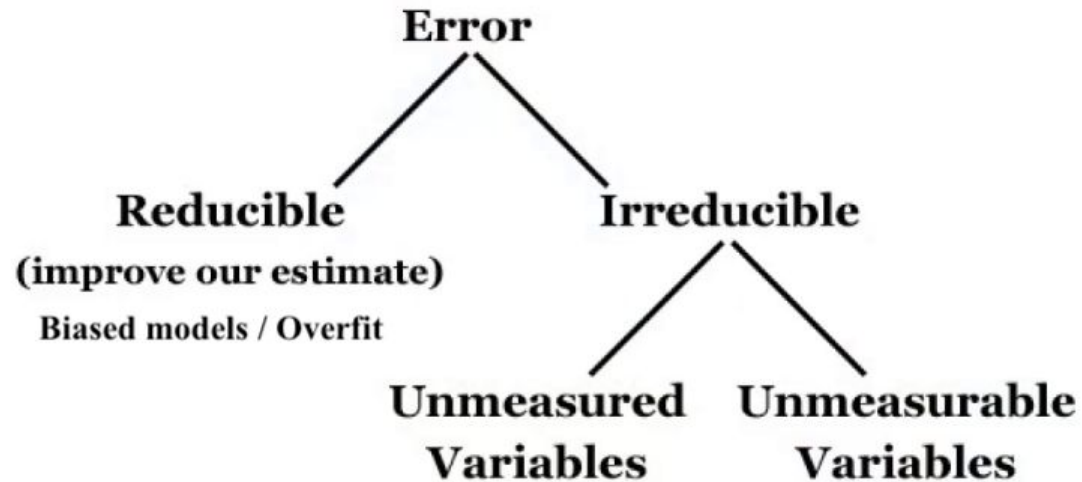
- **.transform():**

This step applies the previously learned transformation to new data (such as the validation or test set) without recalculating any parameters. Used in: Train/Val/Test sets to ensure the same transformation process is applied consistently.

- **.fit_transform():**

Combined operation to save time when you want to learn and apply transformations simultaneously. **ONLY** using on the train dataset.

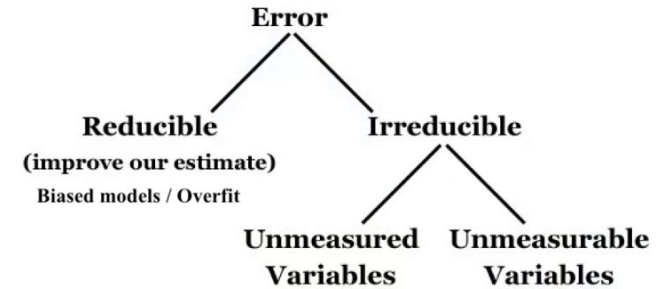
Sources of prediction error



Sources of prediction error

We are given (X, y) training data and we fit a model $f(X)$

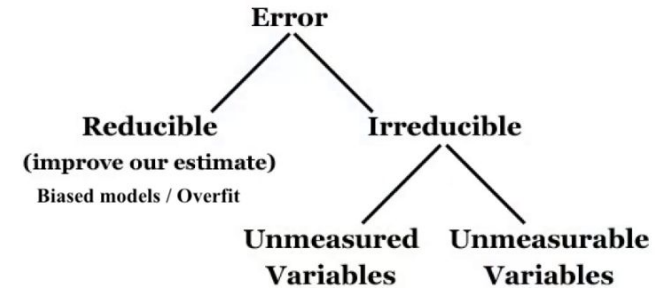
$Err = (f(X_i) - y_i)^2$ from a single observation in the test case



Sources of prediction error

We are given (\mathbf{X}, \mathbf{y}) training data and we fit a model $f(\mathbf{X})$

$Err = (f(\mathbf{X}_i) - y_i)^2$ from a single observation in the test case



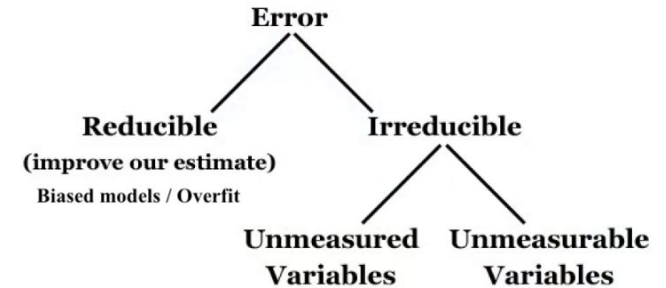
There are three sources of errors in that **Err** number:

1. Noisy \mathbf{X} or \mathbf{y} data, such as inconsistent $\mathbf{X} \rightarrow \mathbf{y}$
2. Model underfitting or bias: Too weak or simple
3. Model overfitting: Model too specific to training data

Sources of prediction error

We are given (\mathbf{X}, \mathbf{y}) training data and we fit a model $f(\mathbf{X})$

$Err = (f(\mathbf{X}_i) - y_i)^2$ from a single observation in the test case



There are three sources of errors in that **Err** number:

1. Noisy \mathbf{X} or \mathbf{y} data, such as inconsistent $\mathbf{X} \rightarrow \mathbf{y}$
2. Model underfitting or bias: Too weak or simple
3. Model overfitting: Model too specific to training data

Conceptually: $Err = \text{"noise"} + \text{"bias"} + \text{"overfitting"}$

Stats nerds: $Err = \text{Irreducible Error} + \text{Bias}^2 + \text{Variance}$

Sources of prediction error

1. Noise can lead to inconsistent data

- Imagine you have two observations such as:
 $[18, 1, 9] \rightarrow 91$
 $[18, 1, 9] \rightarrow 99$
- No model can predict two different y values for the same x vector
- Model will have $Err > 0$ no matter what
- We know this as the **irreducible error**
- Noise comes from faulty sensors, typos, self-reporting issues, etc...
- 🙅 Nothing we can do about the irreducible error

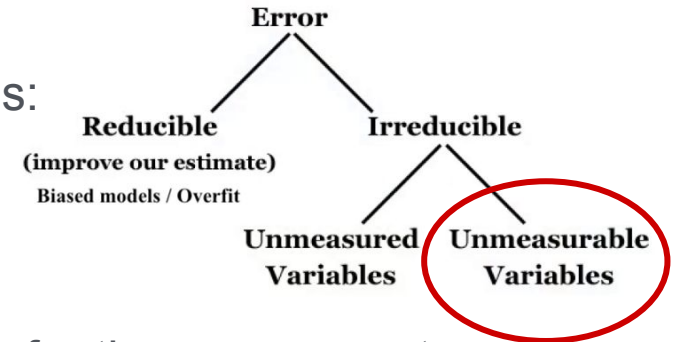
Sources of prediction error

1. Noise can lead to inconsistent data

- Imagine you have two observations such as:

$[18, 1, 9] \rightarrow 91$

$[18, 1, 9] \rightarrow 99$



- No model can predict two different y values for the same x vector
- Model will have $Err > 0$ no matter what
- We know this as the **irreducible error**
- Noise comes from faulty sensors, typos, self-reporting issues, etc...
- 🙅 Nothing we can do about the irreducible error

Sources of prediction error

1. Noise can lead to inconsistent data

- What if inconsistent training observations, such as:
 $[18, 1, 9] \rightarrow 91$
 $[18, 1, 9] \rightarrow 99$
- were really just missing a variable we don't have?
 $[18, 1, 9, 10] \rightarrow 91$
 $[18, 1, 9, 7] \rightarrow 99$
- Example: Two apartment observations look identical: 2bd & 1bath, but they have very different price only because we lack "awesome views" vars
- Missing variables are called *exogenous* variables

Sources of prediction error

1. Noise can lead to inconsistent data

- What if inconsistent training observations, :

[18,1,9] → 91

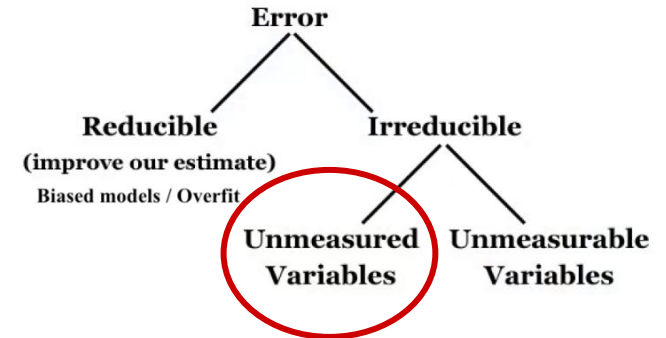
[18,1,9] → 99

- were really just missing a variable we don't have?

[18,1,9,10] → 91

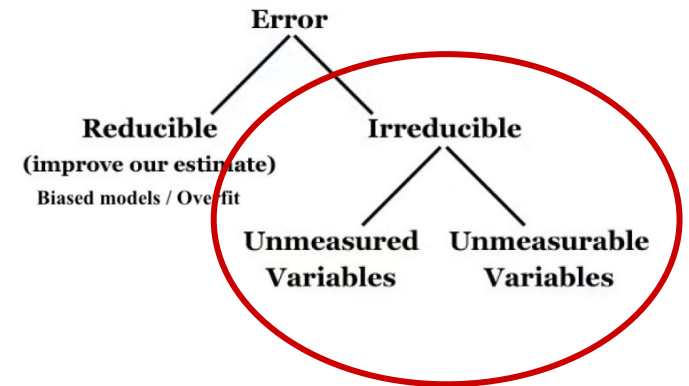
[18,1,9,7] → 99

- Example: Two apartment observations look identical: 2bd & 1bath, but they have very different price only because we lack "awesome views" vars
- Missing variables are called *exogenous* variables

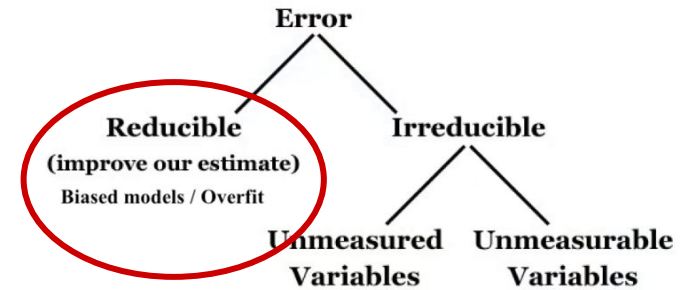


Sources of prediction error

1. Noise can lead to inconsistent data



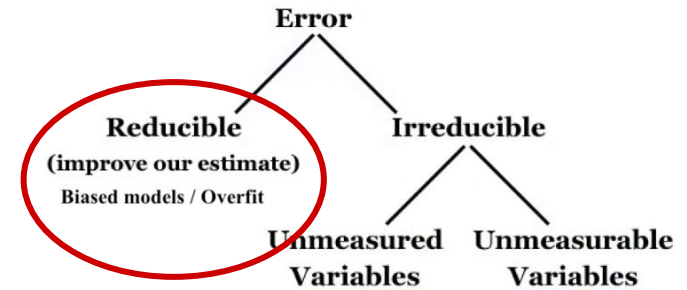
Sources of prediction error



There are three sources of errors in that ***Err*** number:

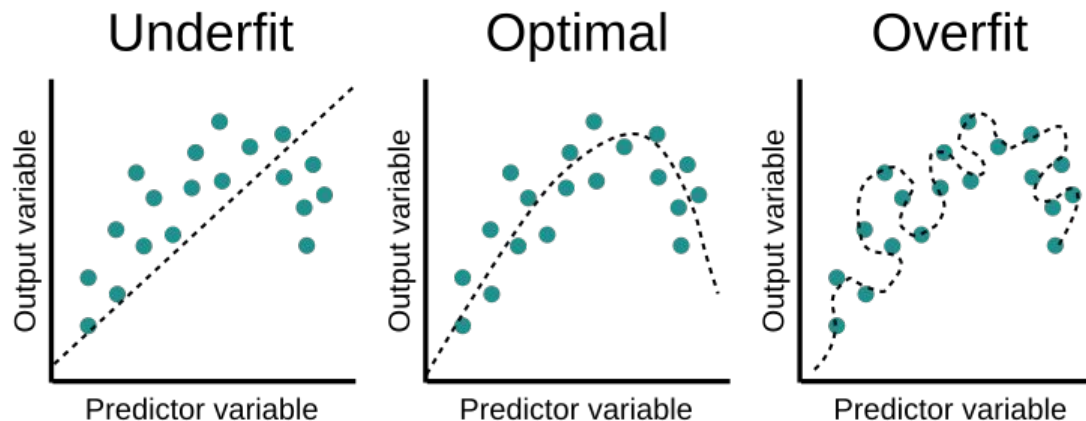
1. Noisy **X** or **y** data, such as inconsistent **X**→**y**
2. Model underfitting or bias: Too weak or simple
3. Model overfitting: Model too specific to training data

Sources of prediction error

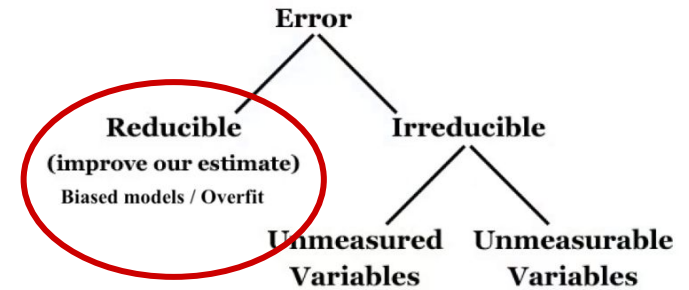


There are three sources of errors in that **Err** number:

1. Noisy **X** or **y** data, such as inconsistent **X**→**y**
2. Model underfitting or bias: Too weak or simple
3. Model overfitting: Model too specific to training data



Sources of prediction error



There are three sources of errors in that **Err** number:

1. Noisy **X** or **y** data, such as inconsistent **X**→**y**
2. Model underfitting or bias: Too weak or simple
3. Model overfitting: Model too specific to training data

Underfit



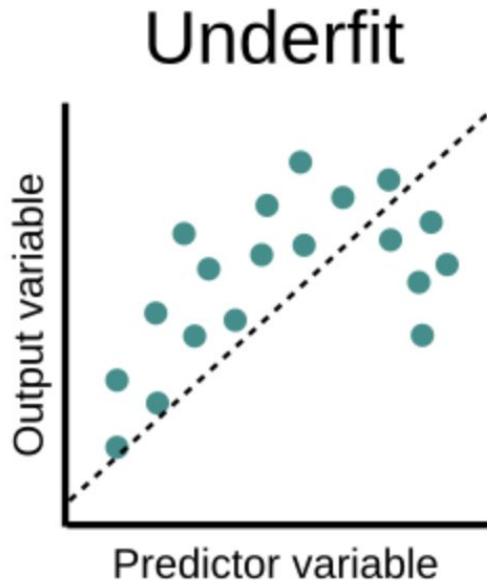
Optimal



Overfit



Overly simple models lead to biased models



- **Bias** is the error rate of your model on the training set
- Bias is how much your model **underfits** the training data

Overly simple models lead to biased models

How do you compute bias?

$$\textit{Bias}[\hat{f}(x)] = E[\hat{f}(x) - y]$$

Expected difference between predicted and observed

Overly simple models lead to biased models



Check-in

A model that has a good ability to fit the training data has _____ bias

Overly simple models lead to biased models

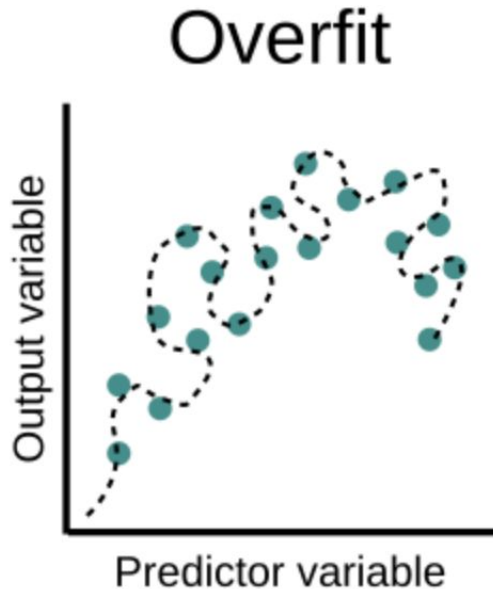


Check-in

A model that has a good ability to fit the training data has **LOW** bias

- We want to minimize bias
- Models with high bias:
 - Fail to capture meaningful patterns in data
 - Under-fit training data
- How to decrease bias? Make model more complex!
 - Add more parameters
 - Pick a different model

Overly complex models can overfit



- **Variance** is the amount a model's prediction will change if a different training data is used, ie, small changes in the training data can result in large changes in the estimated model
- **Variance** in a model is the flexibility to learn patterns in the observed data
- Variance is how much your model **overfits** the training data

Overly complex models can overfit

How do you compute variance?

$$\text{Var}[\hat{f}(x)] = E[\hat{f}(x)^2] - (E[\hat{f}(x)])^2$$

Intuitively, how much the algorithm will move around its mean

Overly complex models can overfit



Check-in

A model that is strongly influenced by the specifics of the training data has _____ variance

Overly complex models can overfit



Check-in

A model that is strongly influenced by the specifics of the training data has **HIGH** variance

- We want to minimize variance. A model that has a good ability to predict test data has **low** variance.
- The more complex the model is, the more data points it will "capture". However, complexity will make the model "move" more to "capture" the data points, and hence its variance will be larger
- How to decrease variance? Make model less complex!
 - A larger training set tends to decrease variance, ie, reduces the chance of overfitting --> increases the chance of generalization
 - Regularization

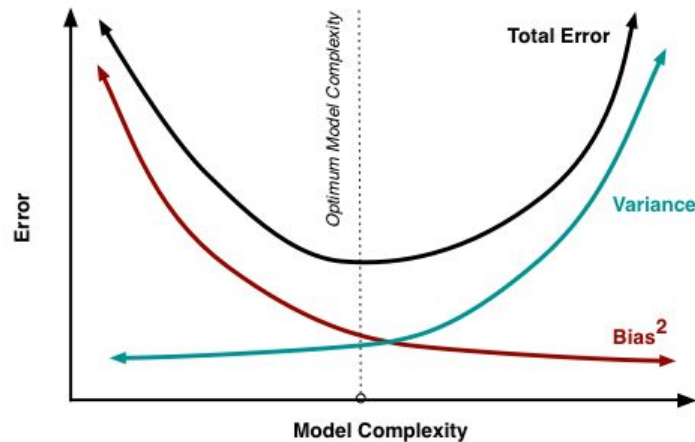
Bias-variance trade-off

So... we need to decrease both bias and variance (to avoid underfitting and overfitting)



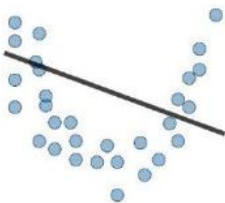


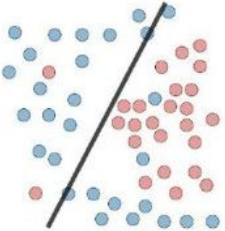
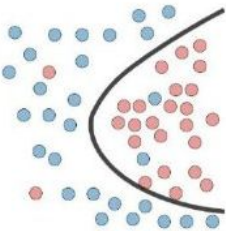
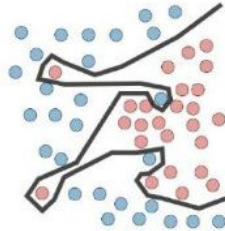

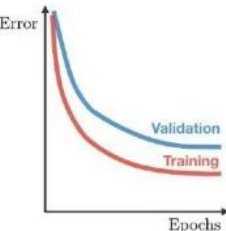
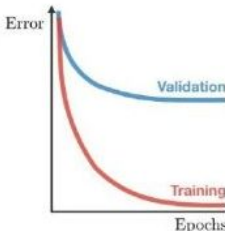


Bias-variance tradeoff



Let the validation dataset be your guide to approximate complexity (typically using cross-validation)

- Collect a lot of data
- Engineer good features
- Pick a complex algorithm
- Train the specific model until validation scores starts to go down (smart early stopping)

	Underfitting	Just right	Overfitting
Symptoms	<ul style="list-style-type: none"> - High training error - Training error close to test error - High bias 	<ul style="list-style-type: none"> - Training error slightly lower than test error 	<ul style="list-style-type: none"> - Low training error - Training error much lower than test error - High variance
Regression			
Classification			
Deep learning			
Remedies	<ul style="list-style-type: none"> - Complexify model - Add more features - Train longer 		<ul style="list-style-type: none"> - Regularize - Get more data

Summary

- **Bias** is the error rate of your model on the training set. Bias is how much your model **underfits** the training data
- **Variance** is the amount a model's prediction will change if a different training data is used. Variance is how much your model **overfits** the training data
- We want to always **minimize both bias and variance**, but when one goes down, the other one goes up (more/less complex model)– hence **bias-variance tradeoff**





Is my model any good?

Performance metrics

- Measurements used to evaluate how well a machine learning model is performing.
- The choice of performance metrics depends on the nature of the problem, whether it's a classification, regression, clustering task, and so on.

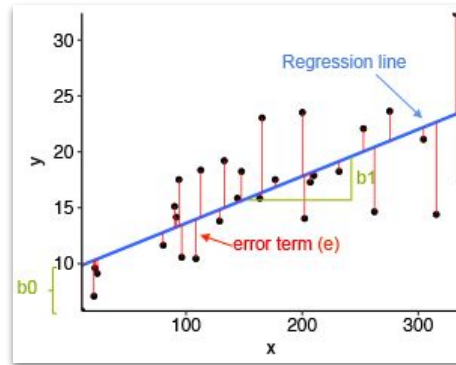
Metrics for regression models	Metrics for classification models
<p>In a regression problem, you're trying to predict a continuous outcome variable (like the price of a house).</p> <p>Performance metrics in regression evaluate the difference between the true and predicted values, often based on errors.</p>	<p>In a classification problem, you're trying to predict which category or class an observation belongs to (like spam vs. non-spam emails).</p> <p>Performance metrics in classification evaluate how well the model can correctly classify the observations.</p>

Performance metrics

- Measurements used to evaluate how well a machine learning model is performing.
- The choice of performance metrics depends on the nature of the problem, whether it's a classification, regression, clustering task, and so on.

Metrics for regression models	Metrics for classification models
<ul style="list-style-type: none">• MAE (Mean Absolute Error)• Average Error• MAPE (Mean Absolute Percentage Error)• RMSE (Root Mean Squared Error)• SST (Total Sum of Squared Errors)• R-Squared• Adjusted R-Squared• and many more... (External link)	<ul style="list-style-type: none">• Misclassification rate• Accuracy• Sensitivity• Specificity• Recall• Precision• F1• ROC-AUC• and many more... (Wiki link)

Performance metrics for regression models



Performance metrics in regression evaluate the difference between the true and predicted values, often based on errors.

MAE (Mean Absolute Error)

- MAE measures the average magnitude of the errors in a set of predictions, without considering their direction.
- It's the average over the test sample of the absolute differences between prediction and actual observation where all individual differences have equal weight.

The diagram illustrates the Mean Absolute Error (MAE) formula with the following components and annotations:

- Formula:** $MAE = \frac{1}{n} \sum |y - \hat{y}|$
- Annotations:**
 - A blue line points from the text "Divide by the total number of data points" to the fraction $\frac{1}{n}$.
 - A green line points from the text "Actual output value" to the variable y inside a green box.
 - A yellow line points from the text "Predicted output value" to the variable \hat{y} inside an orange box.
 - A bracket under the absolute value term $|y - \hat{y}|$ is labeled "The absolute value of the residual".
 - The summation symbol \sum is labeled "Sum of".

MAE (Mean Absolute Error)

- MAE measures the average magnitude of the errors in a set of predictions, without considering their direction.
- It's the average over the test sample of the absolute differences between prediction and actual observation where all individual differences have equal weight.

$$MAE = \frac{1}{n} \sum \left| y - \hat{y} \right|$$

Diagram illustrating the MAE formula components:

- $\frac{1}{n}$: Divide by the total number of data points
- \sum : Sum of
- y : Actual output value
- \hat{y} : Predicted output value
- $|y - \hat{y}|$: The absolute value of the residual

**The lower
the better**

MSE (Mean Square Error)

- Both the mean squared error and the mean absolute error tell you how close a regression line is to a set of points.
- MSE is just like the MAE but squares the difference before summing them all instead of using the absolute value, therefore, it assigns more weight to the bigger errors.

The diagram illustrates the Mean Squared Error (MSE) formula with several annotations:

- A blue box around $\frac{1}{n}$ is labeled "Divide by the total number of data points".
- A green box around the actual output value y is labeled "Actual output value".
- An orange box around the predicted output value \hat{y} is labeled "Predicted output value".
- A bracket under the difference $y - \hat{y}$ is labeled "The square of the difference between actual and predicted".

$$MSE = \frac{1}{n} \sum (y - \hat{y})^2$$

MSE (Mean Square Error)

- Both the mean squared error and the mean absolute error tell you how close a regression line is to a set of points.
- MSE is just like the MAE but squares the difference before summing them all instead of using the absolute value, therefore, it assigns more weight to the bigger errors.

$$MSE = \frac{1}{n} \sum (y - \hat{y})^2$$

Diagram annotations for the MSE formula:

- Divide by the total number of data points**: Points to the $\frac{1}{n}$ term.
- Actual output value**: Points to the y term inside the parentheses.
- Predicted output value**: Points to the \hat{y} term inside the parentheses.
- The square of the difference between actual and predicted**: Points to the entire $(y - \hat{y})^2$ term.

**The lower
the better**

RMSE (Root Mean Square Error)

- It shares MSE's property of heavily penalizing larger errors (because it's based on the squared differences), but it's in the same units as the outcome variable, which can make it easier to interpret.

The diagram illustrates the formula for Root Mean Square Error (RMSE) with hand-drawn annotations. The formula is $\sqrt{MSE} = \sqrt{\frac{1}{n} \sum (y - \hat{y})^2}$. Annotations include: a blue box around $\frac{1}{n}$ with the text "Divide by the total number of data points"; a green box around y with the text "Actual output value"; an orange box around \hat{y} with the text "Predicted output value"; and a bracket under the difference $y - \hat{y}$ with the text "The square of the difference between actual and predicted".

$$\sqrt{MSE} = \sqrt{\frac{1}{n} \sum (y - \hat{y})^2}$$

Divide by the total number of data points

Actual output value

Predicted output value

The square of the difference between actual and predicted

RMSE (Root Mean Square Error)

- It shares MSE's property of heavily penalizing larger errors (because it's based on the squared differences), but it's in the same units as the outcome variable, which can make it easier to interpret.

The diagram illustrates the formula for Root Mean Square Error (RMSE). The formula is written as \sqrt{MSE} , where $MSE = \frac{1}{n} \sum (y - \hat{y})^2$. Annotations include: a blue box around $\frac{1}{n}$ with the text "Divide by the total number of data points"; a green box around y with the text "Actual output value"; an orange box around \hat{y} with the text "Predicted output value"; and a bracket under the difference $y - \hat{y}$ with the text "The square of the difference between actual and predicted".

$$\sqrt{MSE} = \sqrt{\frac{1}{n} \sum (y - \hat{y})^2}$$

The lower
the better

R-Squared

- R-squared, also known as the coefficient of determination, measures the proportion of the variance in the dependent variable that can be predicted from the independent variable(s).
- In other words, it shows how close the data are to the fitted regression line. An R-squared of 1 indicates that all changes in the dependent variable are completely explained by changes in the independent variable(s).

$$R^2 = 1 - \frac{SS_{RES}}{SS_{TOT}} = 1 - \frac{\sum_i (y_i - \hat{y}_i)^2}{\sum_i (y_i - \bar{y})^2}$$


SSRES: the
residual sum of
squared errors

SSTOT: the total
sum of squared
errors

R-Squared

- R-squared, also known as the coefficient of determination, measures the proportion of the variance in the dependent variable that can be predicted from the independent variable(s).
- In other words, it shows how close the data are to the fitted regression line. An R-squared of 1 indicates that all changes in the dependent variable are completely explained by changes in the independent variable(s).

$$R^2 = 1 - \frac{SS_{RES}}{SS_{TOT}} = 1 - \frac{\sum_i (y_i - \hat{y}_i)^2}{\sum_i (y_i - \bar{y})^2}$$



**The higher
the better**

Adjusted R-Squared

- There's a potential problem with R^2 : it never decreases as more variables are added to the model, even if those variables are only weakly associated with the response. This is where Adjusted R-squared comes in.
- Adjusted R-Squared also shows how well the regression line approximates the real data points but it also takes into account the number of predictors in the model.

$$\text{Adjusted } R^2 = 1 - \frac{(1 - R^2)(N - 1)}{N - p - 1}$$


N = the number of
observations

p = the number of
predictors

Adjusted R-Squared

- There's a potential problem with R^2 : it never decreases as more variables are added to the model, even if those variables are only weakly associated with the response. This is where Adjusted R-squared comes in.
- Adjusted R-Squared also shows how well the regression line approximates the real data points but it also takes into account the number of predictors in the model.

$$\text{Adjusted } R^2 = 1 - \frac{(1 - R^2)(N - 1)}{N - p - 1}$$



**The higher
the better**

Which metric should I get started with?

- Among the many available metrics, these two are good starting points because they provide insights into different aspects:
 - **RMSE** (Root Mean Squared Error) measures the magnitude of the error in real units. (We want this to be low.)
 - **R²** (R-squared) tells us how well the model explains the variability in the data. (We want this to be high, close to 1.)
- Example:
 - If the RMSE is 500 in a model that predicts house prices in thousands of dollars, this means the model has an average prediction error of \$500,000. If the R² is 0.85, the model explains 85% of the variation in house prices.

Summary: performance metrics for regression models

- [Cheat Sheet](#)

Train/val/test data sets



- It is important to split the data into train/val/test **BEFORE** doing any **feature engineering** or modeling to prevent **data leakage** and ensure that the model is truly evaluated on unseen data

Data leakage: When information from val or test is used to inform the model during training or feature engineering



No peeking!!!

The test data set is used **ONLY** after you think you have the best model. It's the only true measure of generality.



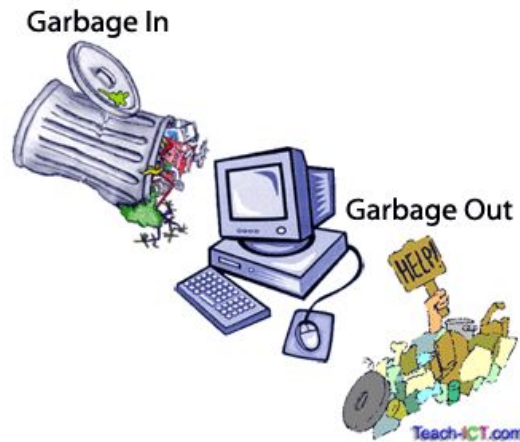
What We See

08 02 22 97 38 15 00 40 00 75 04 05 07 78 52 12 50 77 91 08 08 02 22 97
 49 49 99 40 17 81 18 57 60 87 17 40 98 43 69 48 04 56 62 00 49 49 99 40
 81 49 31 73 55 79 14 29 93 71 40 67 53 88 30 03 49 13 36 65 81 49 31 73
 52 70 95 23 04 60 11 42 69 24 68 56 01 32 56 71 37 02 36 91 52 70 95 23
 22 31 16 71 51 67 63 89 41 92 36 54 22 40 40 28 66 33 13 80 22 31 16 71
 24 47 32 60 99 03 45 02 44 75 33 53 78 36 84 20 35 17 12 50 24 47 32 60
 32 98 81 28 64 23 67 10 26 38 40 67 59 54 70 66 18 38 64 70 32 98 81 28
 67 26 20 68 02 62 12 20 95 63 94 39 63 08 40 91 66 49 94 21 67 26 20 68
 24 55 58 05 66 73 99 26 97 17 78 78 96 83 14 88 34 89 63 72 24 55 58 05
 21 36 23 09 75 00 76 44 20 45 35 14 00 61 33 97 34 31 33 95 21 36 23 09
 78 17 53 28 22 75 31 67 15 94 03 80 04 62 16 14 09 53 56 92 78 17 53 28
 16 39 05 42 96 35 31 47 55 58 88 24 00 17 54 24 36 29 85 57 16 39 05 42
 86 56 00 48 35 71 89 07 05 44 44 37 44 60 21 58 51 54 17 58 86 56 00 48
 19 80 81 68 05 94 47 69 28 73 92 13 86 52 17 77 04 89 55 40 19 80 81 68
 04 52 08 83 97 35 99 16 07 97 57 32 16 26 26 79 33 27 98 66 04 52 08 83
 88 36 68 87 57 62 20 72 03 46 33 67 46 55 12 32 63 93 53 69 88 36 68 87
 04 42 16 73 38 25 39 11 24 94 72 18 08 46 29 32 40 62 76 36 04 42 16 73
 20 69 36 41 72 30 23 88 34 62 99 69 82 67 59 85 74 04 36 16 20 69 36 41
 20 73 35 29 78 31 90 01 74 31 49 71 48 86 81 16 23 57 05 54 20 73 35 29
 01 70 54 71 83 51 54 69 16 92 33 48 61 43 52 01 89 19 67 48 01 70 54 71

What Computers See

Feature Engineering

- Feature engineering is the process of selecting and transforming raw data into features that can be used to train machine learning models
- AKA creating/changing number of columns



Feature Engineering

Handling categorical attributes

Most machine learning models prefer to work with numbers

- **Ordinal encoder**

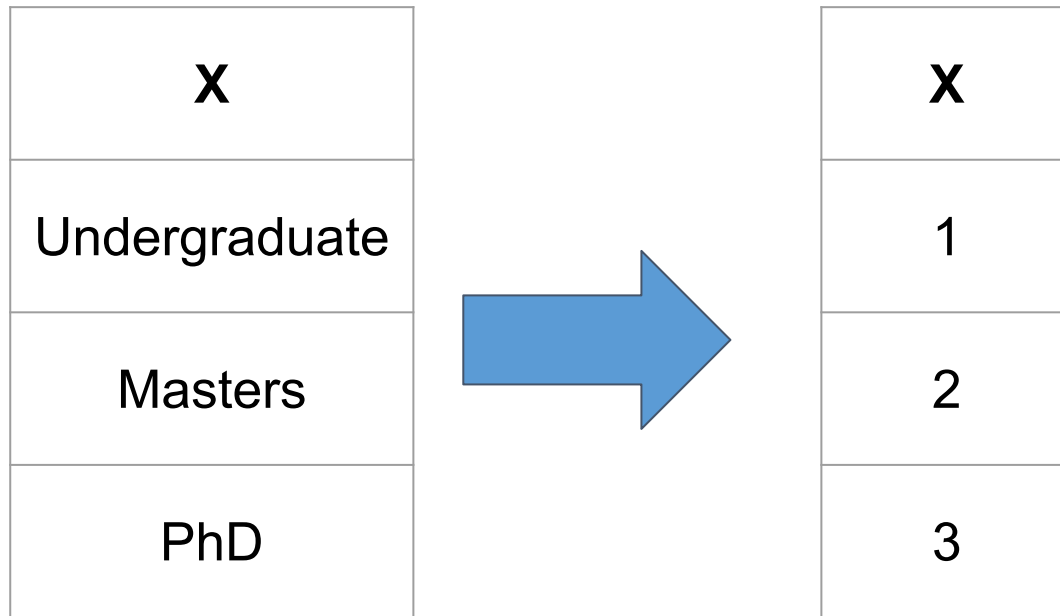
Disadvantage: Assumes two nearby values are more similar to each other (useful for ordered categories such as "bad", "average", "good", "excellent")

- **One-Hot encoder**

One binary category per attribute

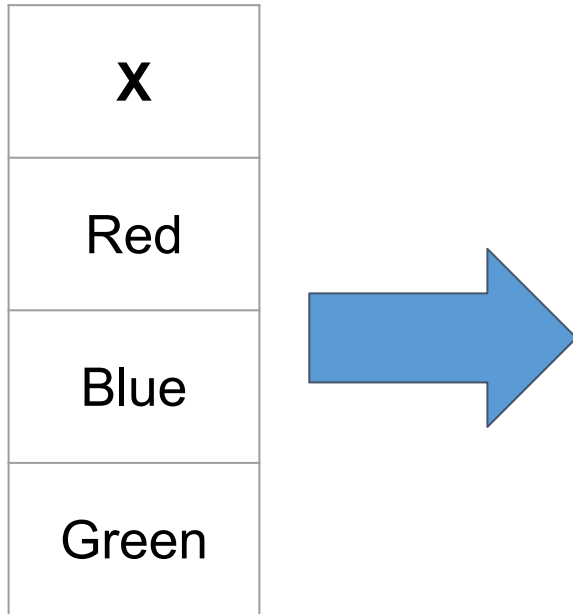
Ordinal encoder

Example: $X = [\text{"Undergraduate"}, \text{"Masters"}, \text{"PhD"}]$



One-Hot Encoding

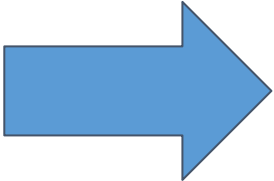
Example: $X = [\text{"Red"}, \text{"Blue"}, \text{"Green"}]$



One-Hot Encoding

Example: $X = [\text{"Red"}, \text{"Blue"}, \text{"Green"}]$

X			
Red	1	0	0
Blue	0	1	0
Green	0	0	1



X_red	X_blue	X_green
1	0	0
0	1	0
0	0	1

Feature Engineering

Feature Scaling

With few exceptions, machine learning algorithms don't perform well when the input of the numerical values have very different scales

Example in this dataset: total_rooms, median_income

- **MinMax**

- AKA normalization
- Values are shifted and rescaled so that they end up ranging from 0 to 1
- How to calculate it? Subtracting the min val and diving by the difference between the min and the max

- **Standardization**

- How to calculate it? Subtracts the mean value, then it divides the result by the standard deviation
- It doesn't restrict the values to a specific range

MinMax

Example: $X = [2, 10, 15]$

$$= \frac{X - X_{\min}}{X_{\max} - X_{\min}}$$

2:

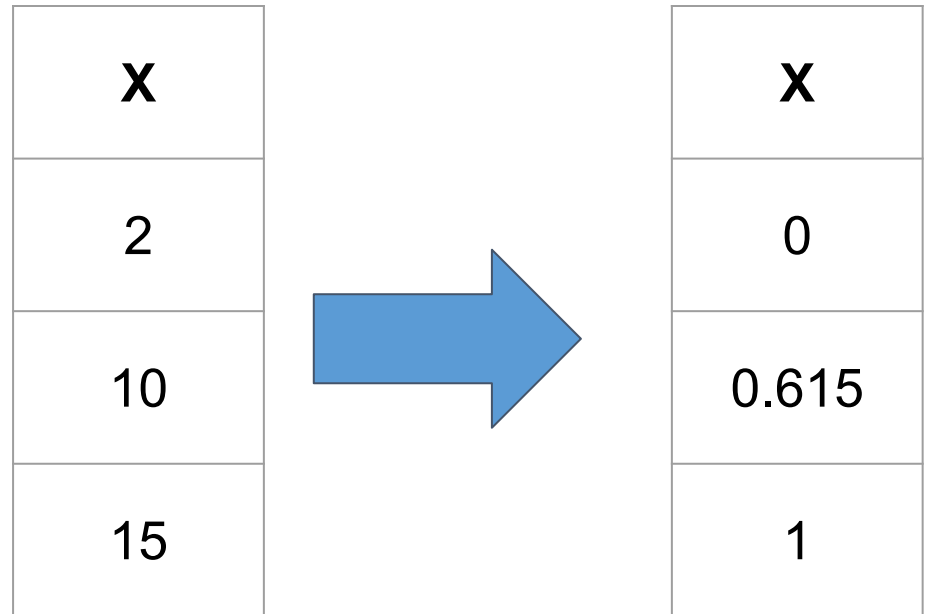
$$\frac{2 - 2}{15 - 2} = \frac{0}{13} = 0$$

10:

$$\frac{10 - 2}{15 - 2} = \frac{8}{13} \approx 0.615$$

15:

$$\frac{15 - 2}{15 - 2} = \frac{13}{13} = 1$$



Standardization

Example: $X = [2, 10, 15]$

- Media (μ):

$$\mu = \frac{2 + 10 + 15}{3} = \frac{27}{3} = 9$$

- Desviación estándar (σ):

$$\sigma = \sqrt{\frac{(2-9)^2 + (10-9)^2 + (15-9)^2}{3}} = \sqrt{\frac{49 + 1 + 36}{3}} = \sqrt{\frac{86}{3}} \approx 5.387$$

$$= \frac{X - \mu}{\sigma}$$

2:

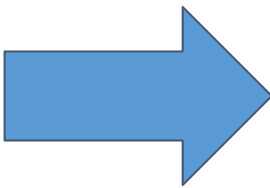
$$\frac{2 - 9}{5.387} = \frac{-7}{5.387} \approx -1.299$$

10:

$$\frac{10 - 9}{5.387} = \frac{1}{5.387} \approx 0.186$$

15:

$$\frac{15 - 9}{5.387} = \frac{6}{5.387} \approx 1.114$$

X		X
2		-1.299
10		0.186
15		1.114

Feature Engineering

Included but not limited to:

- Handling categorical features: Ordinal encoder, OneHot encoder, etc.
- Feature scaling: MinMax, standarizaton, square root, log, etc.
- Feature extraction: Creating new features from existing ones
- Text processing (TF-IDF, Word2Vec, BERT, etc.)
- Extracting things that make sense (date→holiday)
- Dimensionality reduction (PCA, t-SNE, etc.)

Depends on your dataset!

Feature Engineering

Fit vs Fit_Transform vs Transform ⚠

Feature Engineering

Fit vs Fit_Transform vs Transform

FIT

- What it does: The fit method computes the necessary parameters needed to apply a transformation or learn from the data. However, it doesn't actually modify the data.
- When to use it: Use fit when you want to compute the parameters based on your data (like mean, standard deviation for standard scaling, or the min and max values for MinMax scaling), but you don't want to transform the data yet.
- Here, the scaler learns the mean and standard deviation from X_train, but X_train remains unchanged.

```
from sklearn.preprocessing import StandardScaler  
scaler = StandardScaler()  
scaler.fit(X_train)
```

Feature Engineering

Fit vs Fit_Transform vs Transform ⚠

TRANSFORM

- What it does: The transform method applies a previously computed transformation to the data. It doesn't compute any new parameters; it simply uses the parameters computed during the fit phase.
- When to use it: Use transform when you want to apply a transformation to new or existing data using parameters learned during the fit phase.
- Here, both `X_train` and `X_test` are scaled using the mean and standard deviation learned from `X_train` during the fit phase.

```
X_train_scaled = scaler.transform(X_train)
X_test_scaled = scaler.transform(X_test)
```

Feature Engineering

Fit vs Fit_Transform vs Transform ⚠

FIT_TRANSFORM

- What it does: `fit_transform` is essentially a combination of `fit` and `transform`. It computes the necessary parameters from the data and then immediately applies the transformation.
- When to use it: Use `fit_transform` when you're sure you want to both compute the parameters and transform your data in one step. This can be more efficient than calling `fit` and `transform` separately.
- This achieves the same result as the two-step process above but in a single step.

```
x_train_scaled = scaler.fit_transform(X_train)
```

Feature Engineering

Fit vs Fit_Transform vs Transform ⚠

IMPORTANT:

When splitting data into training and test sets, **always fit (or fit_transform) on the training data only**. This simulates a real-world scenario where your model has to generalize to unseen data. By fitting only on the training data, you avoid leaking information from the test set into the training process.

For instance, when scaling:

1. Use fit_transform (or fit and transform) on the training data.
2. Use transform on the test data (do not re-fit the scaler on test data).

This ensures that the test data is transformed based on the parameters learned from the training data, preserving the integrity of the evaluation process.

Example

`X = [1, 2, 3, 4, 5, 6, 7, 8, 9, 10]`

- Let's randomly split this dataset into train/test

`X_train = [1, 2, 3, 6, 7, 8]`

`X_test = [4, 5, 9, 10]`

Example

`X = [1, 2, 3, 4, 5, 6, 7, 8, 9, 10]`

- Let's randomly split this dataset into train/test

`X_train = [1, 2, 3, 6, 7, 8]`

`X_test = [4, 5, 9, 10]`

1. First step: Let's fit the scaler on the train set:

`MinMax.fit(X_train)`

Minimum value = 1

Maximum = 8

Example

$X = [1, 2, 3, 4, 5, 6, 7, 8, 9, 10]$

- Let's randomly split this dataset into train/test

$X_{\text{train}} = [1, 2, 3, 6, 7, 8]$

$X_{\text{test}} = [4, 5, 9, 10]$

- First step: Let's fit the scaler on the train set:

`MinMax.fit(X_train)`

Minimum value = 1

Maximum = 8

- Second step: Let's transform our train data

$X_{\text{train}} = 1:$

$$\frac{1-1}{8-1} = \frac{0}{7} = 0$$

$X_{\text{train}} = 2:$

$$\frac{2-1}{8-1} = \frac{1}{7} \approx 0.143$$

$X_{\text{train}} = 3:$

$$\frac{3-1}{8-1} = \frac{2}{7} \approx 0.286$$

$X_{\text{train}} = 6:$

$$\frac{6-1}{8-1} = \frac{5}{7} \approx 0.714$$

$X_{\text{train}} = 7:$

$$\frac{7-1}{8-1} = \frac{6}{7} \approx 0.857$$

$X_{\text{train}} = 8:$

$$\frac{8-1}{8-1} = \frac{7}{7} = 1$$

Example

$X = [1, 2, 3, 4, 5, 6, 7, 8, 9, 10]$

- Let's randomly split this dataset into train/test

$X_{\text{train}} = [1, 2, 3, 6, 7, 8]$

$X_{\text{test}} = [4, 5, 9, 10]$

- First step: Let's fit the scaler on the train set:

`MinMax.fit(X_train)`

Minimum value = 1

Maximum = 8

- Second step: Let's transform our train data

- Third step: Let's transform our test data

$X_{\text{train}} = 1:$

$$\frac{1-1}{8-1} = \frac{0}{7} = 0$$

$X_{\text{train}} = 2:$

$$\frac{2-1}{8-1} = \frac{1}{7} \approx 0.143$$

$X_{\text{train}} = 3:$

$$\frac{3-1}{8-1} = \frac{2}{7} \approx 0.286$$

$X_{\text{train}} = 6:$

$$\frac{6-1}{8-1} = \frac{5}{7} \approx 0.714$$

$X_{\text{train}} = 7:$

$$\frac{7-1}{8-1} = \frac{6}{7} \approx 0.857$$

$X_{\text{train}} = 8:$

$$\frac{8-1}{8-1} = \frac{7}{7} = 1$$

$X_{\text{test}} = 4:$

$$\frac{4-1}{8-1} = \frac{3}{7} \approx 0.429$$

$X_{\text{test}} = 5:$

$$\frac{5-1}{8-1} = \frac{4}{7} \approx 0.571$$



$X_{\text{test}} = 9:$

$$\frac{9-1}{8-1} = \frac{8}{7} \approx 1.143$$

$X_{\text{test}} = 10:$

$$\frac{10-1}{8-1} = \frac{9}{7} \approx 1.286$$




AGENDA

-  Required activities for Module 8
-  Content review Module 8: Feature Engineering and Overfitting
- **Code**
- Questions

Code

- https://drive.google.com/file/d/1ccaTiK4kli_WMikBP_Q20YPWgCUmkXwF/view?usp=sharing





AGENDA

-  Required activities for Module 8
-  Content review Module 8: Feature Engineering and Overfitting
-  Code
- Questions

QUESTIONS?



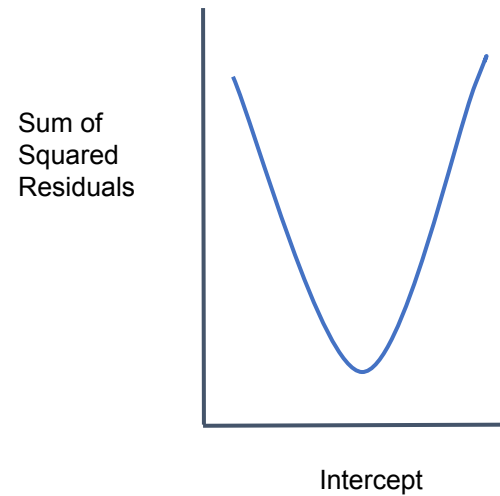
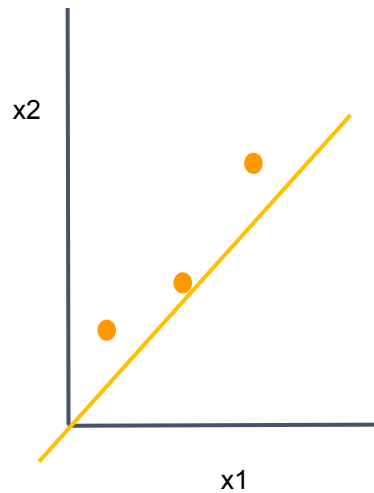
AGENDA

-  Required activities for Module 8
-  Content review Module 8: Feature Engineering and Overfitting
-  Code
-  Questions

APPENDIX

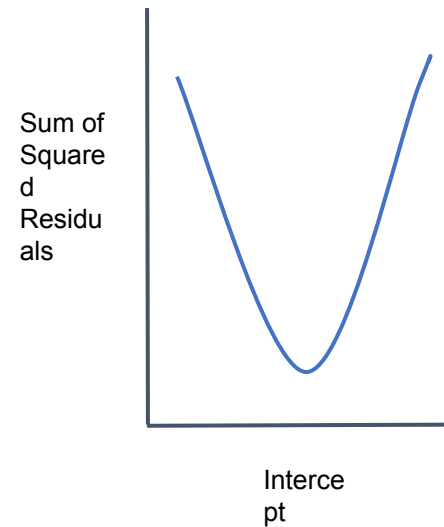
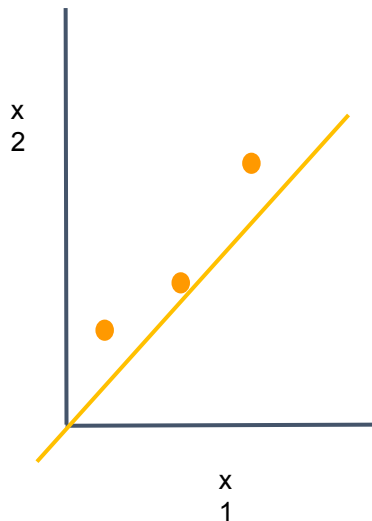
Objective of a Neural Network:

Find weights and biases that minimizes the loss function



Objective of a Neural Network:

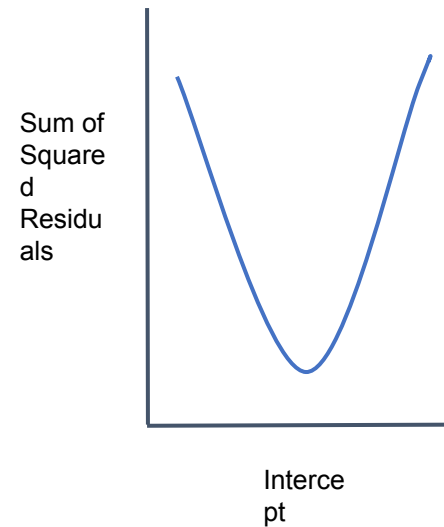
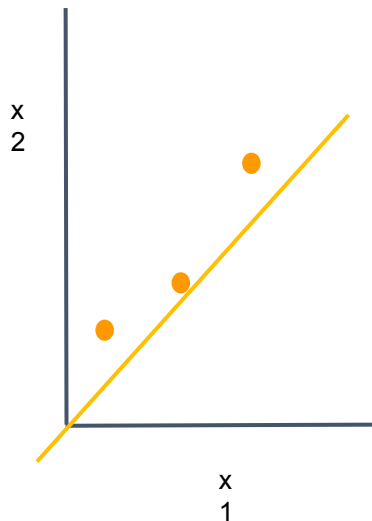
Find weights and biases that minimizes the loss function



Derivative: Slope of a curve at a specific point

Objective of a Neural Network:

Find weights and biases that minimizes the loss function



Derivative: Slope of a curve at a specific point

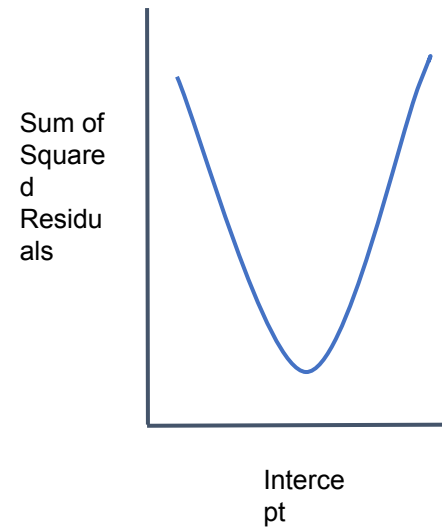
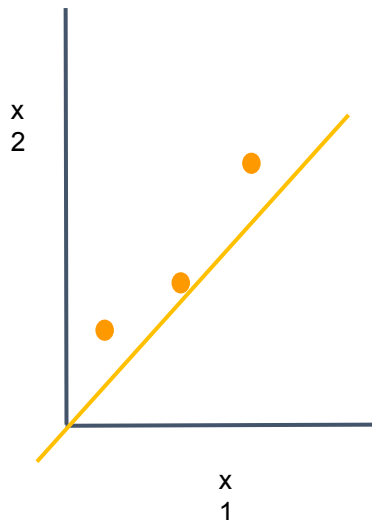
Gradient: Generalization of the derivative. It extends this idea to function of multiple variables.

Objective of a Neural Network:

Find weights and biases that minimizes the loss function



A gradient is like an arrow that points in the direction of the steepest increase on a surface.



If you imagine being on a hill and you want to know which way is uphill, the gradient would point in that direction. If you want to go downhill (like in gradient descent), you'd go in the opposite direction of the gradient.

Derivative: Slope of a curve at a specific point

Gradient: Generalization of the derivative. It extends this idea to function of multiple variables.

Gradient descent

Objective of a Neural Network:

Find weights and biases that minimizes the loss function

- The **gradient** of the loss gives you the direction of the steepest ascent
- To minimize, you move the opposite way (towards the steepest descent)

