

# Section 2 Wrap-Up

Codestates AIB  
JHLEE

# Sprint 1 Linear Models

Supervised learning

Regression

Regularization(Ridge, Lasso)

Classification

Baseline model

One-hot encoding

Train/Validation/Test sets

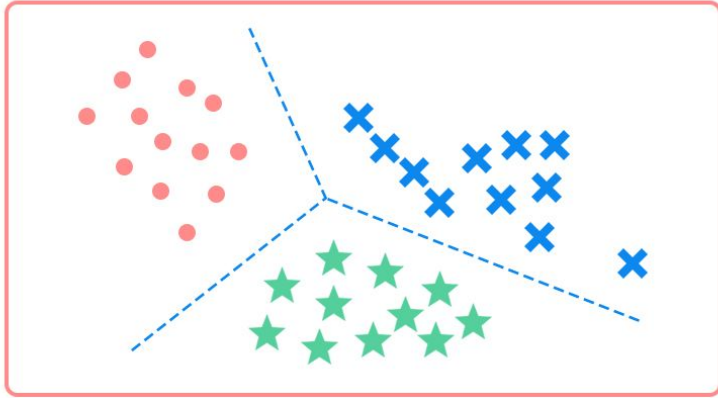
Overfitting/Underfitting

# Supervised Learning



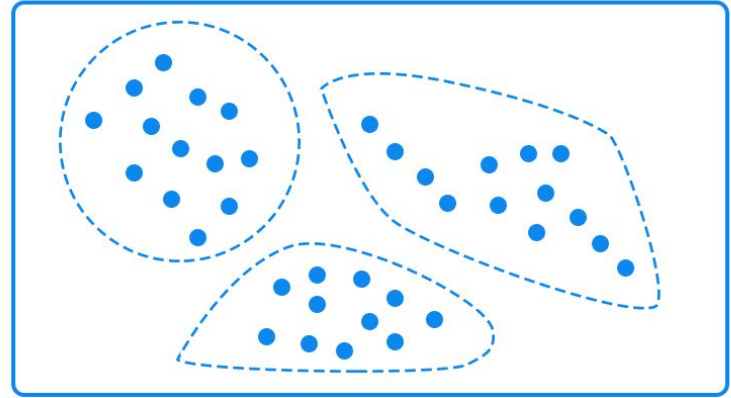
## Supervised vs. Unsupervised Learning

Classification



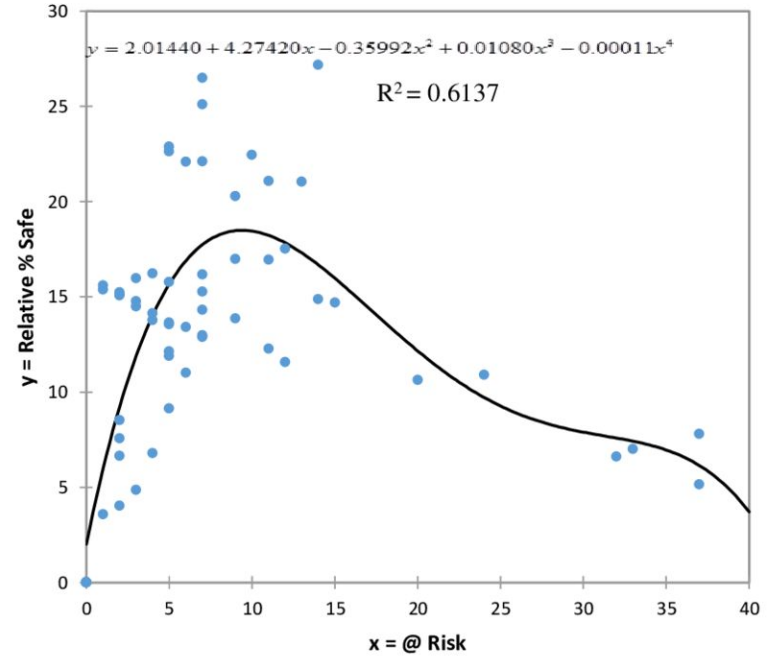
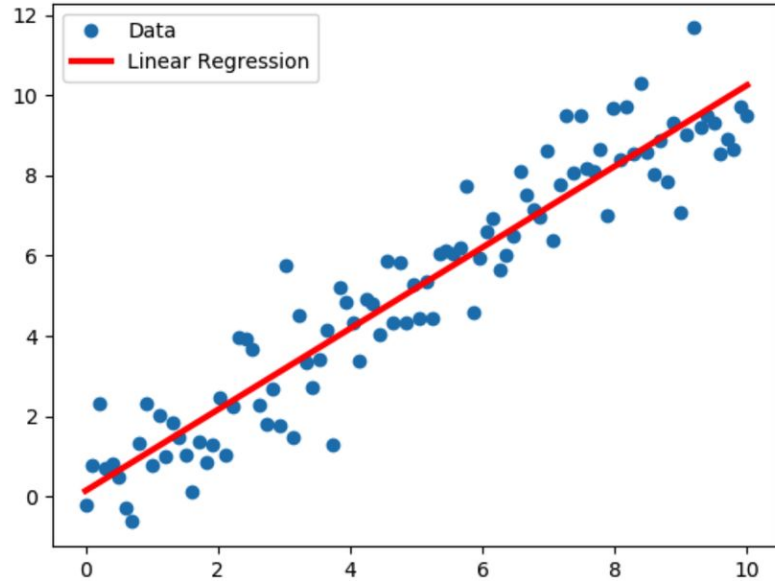
Supervised learning

Clustering

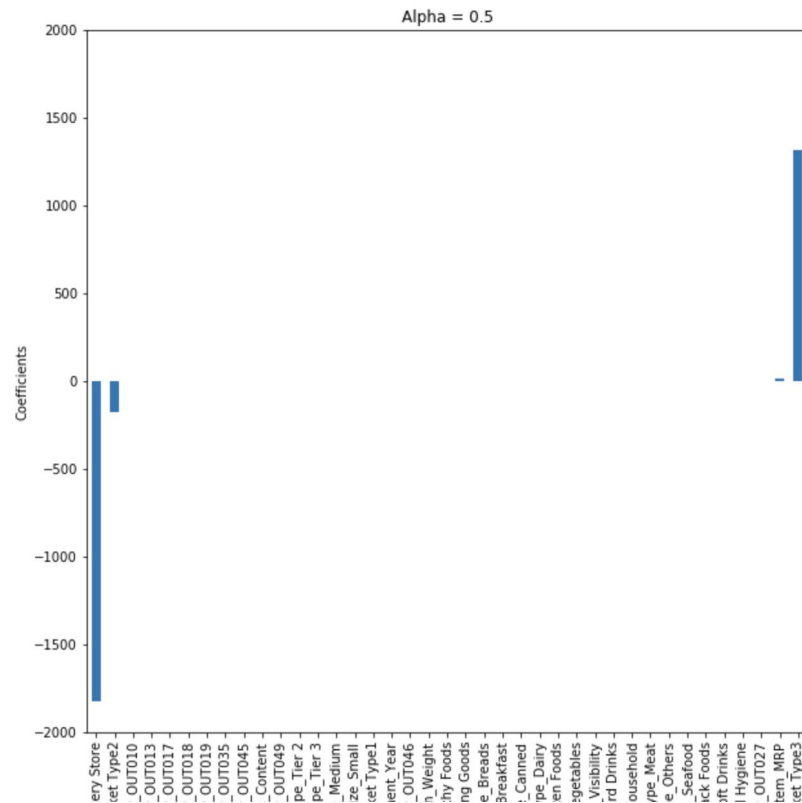
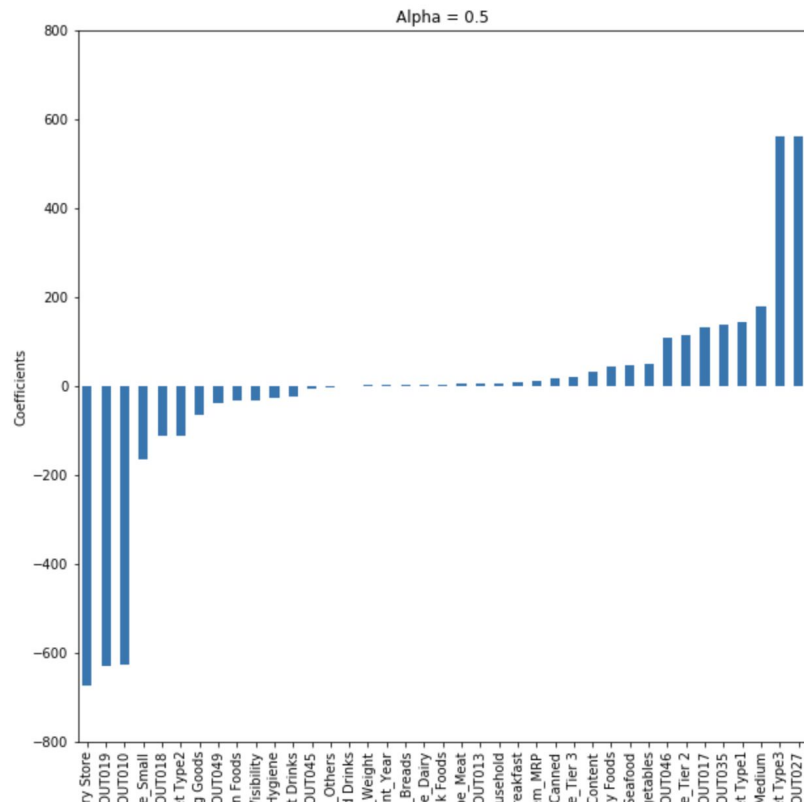


Unsupervised learning

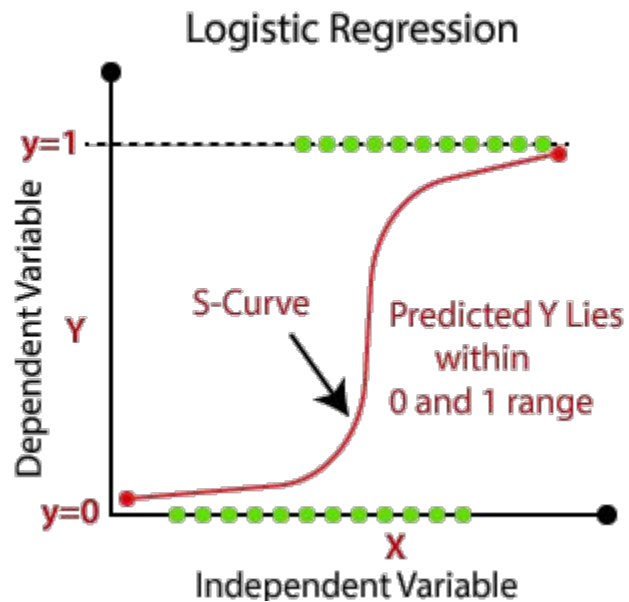
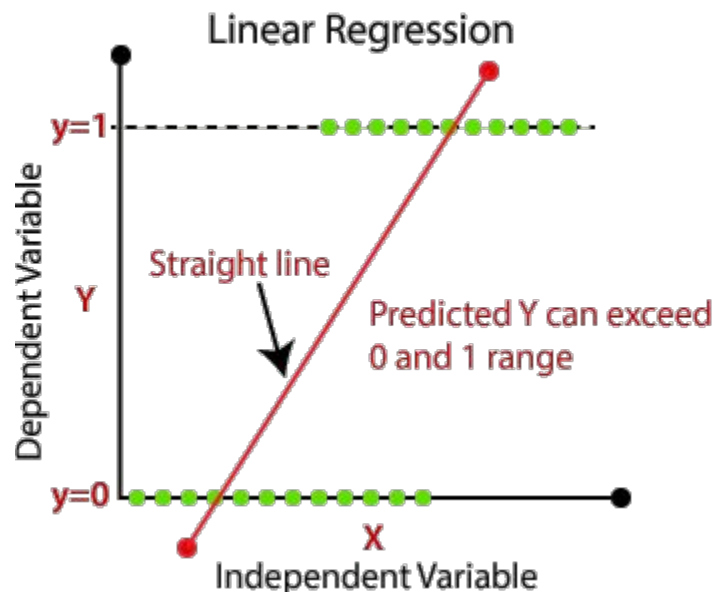
# Regression



# Regularization(Ridge, Lasso)

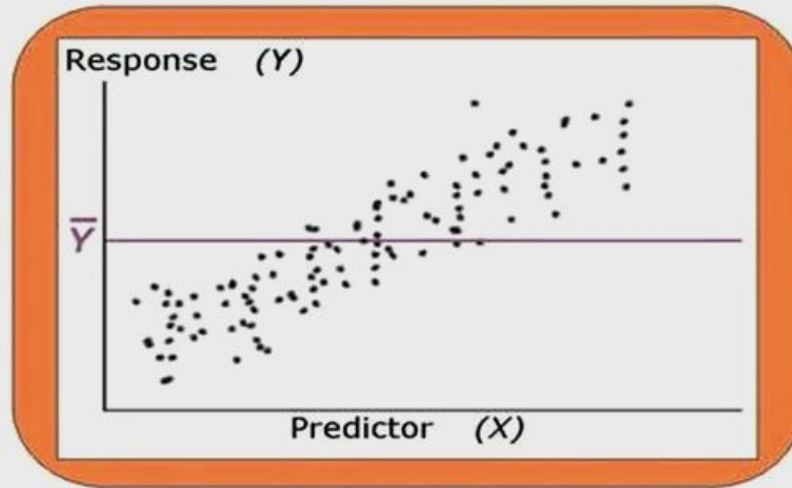


# Classification



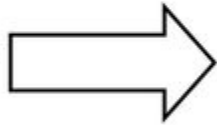
# Baseline model

## The Baseline Model



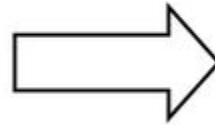
# One-hot encoding

**"Red sticky clay"**



**Numerical**

Word	ID
Red	0
Sticky	1
Clay	2



**One-hot**

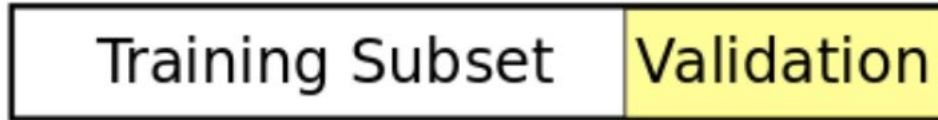
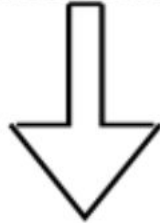
Red	Sticky	Clay
1	0	0
0	1	0
0	0	1



# Train/Validation/Test sets



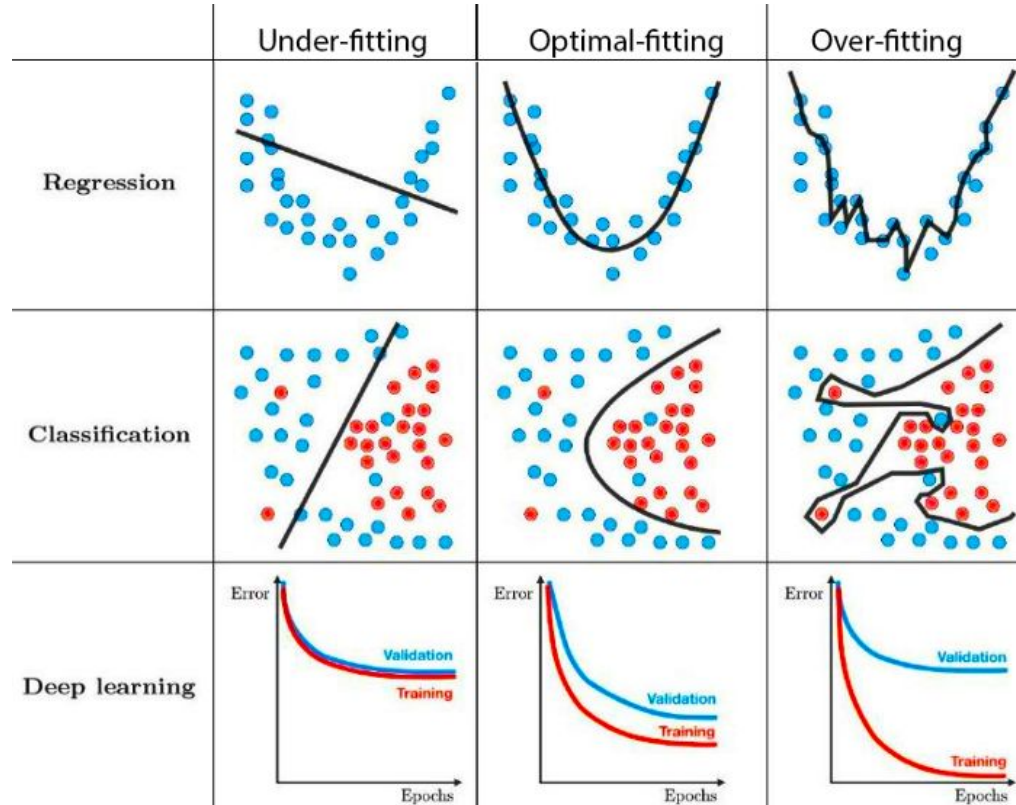
Get Test set error



CV Loop

Tune hyperparameters

# Overfitting/Underfitting



# Sprint 2 Tree Based Models

Decision tree

Sklearn pipeline

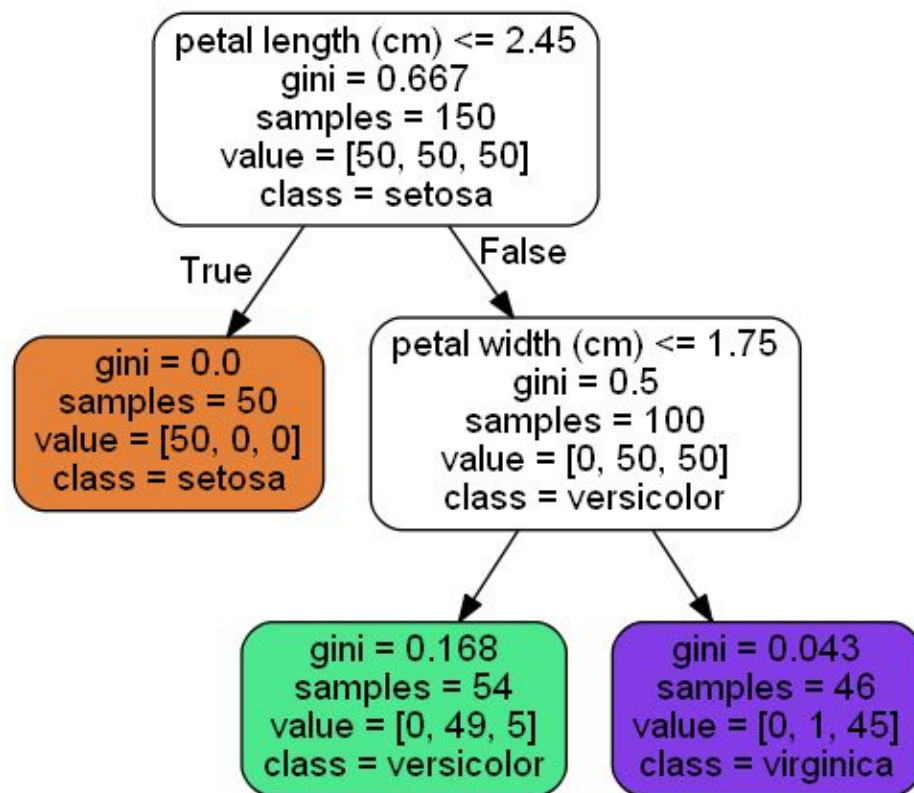
Ensemble(Random forests)

Classification metrics(Accuracy, Precision, Recall, F1, ROC, AUC)

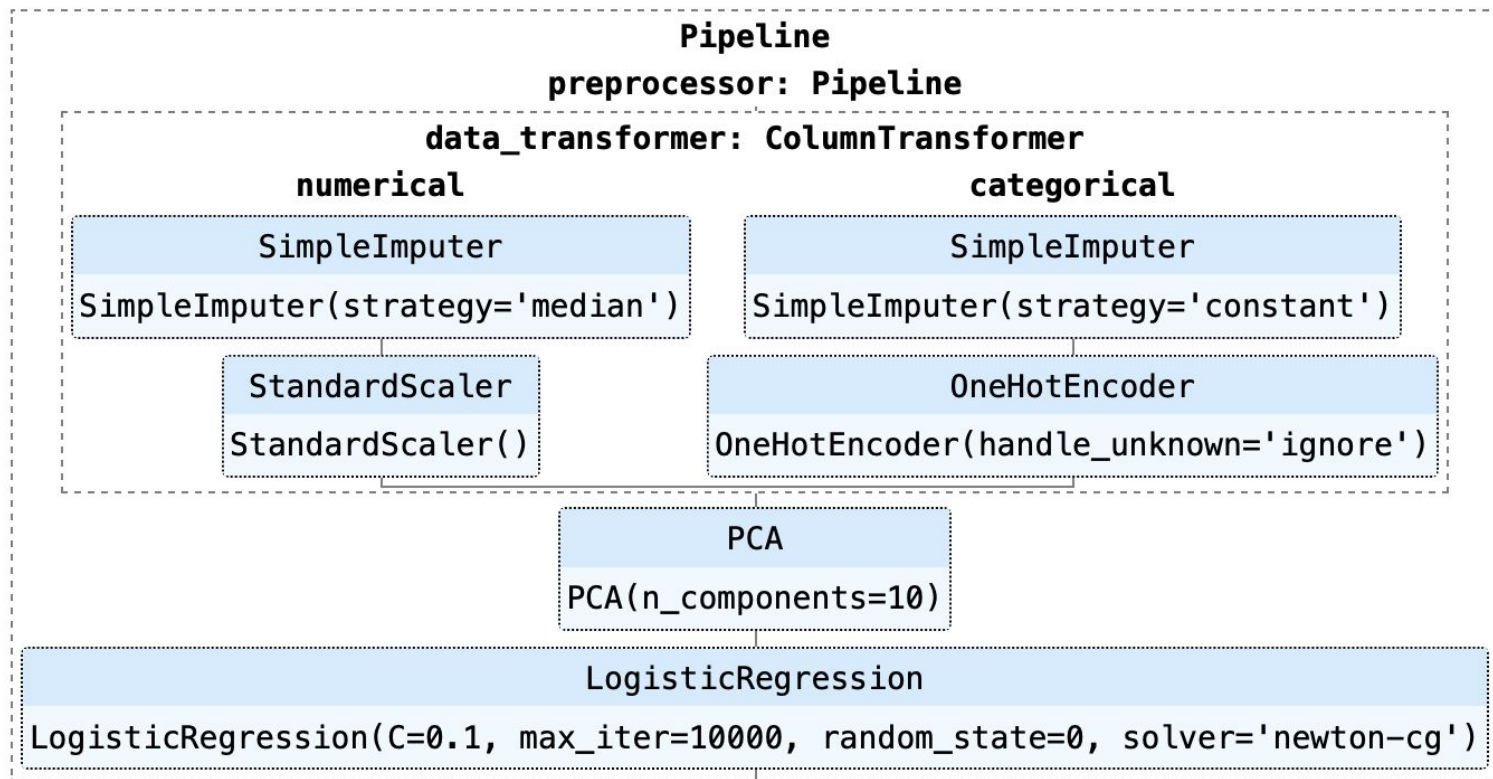
Cross-validation

Hyperparameter tuning

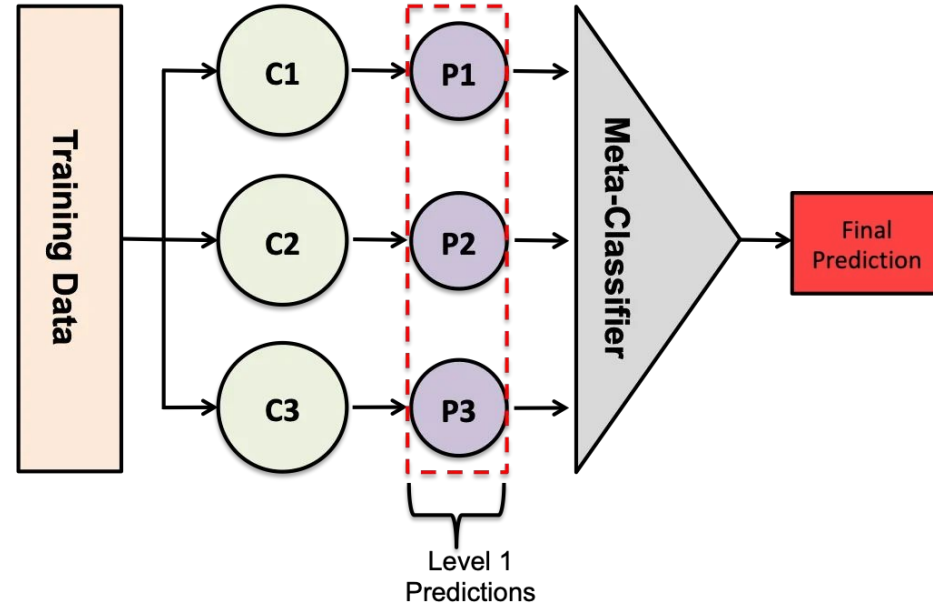
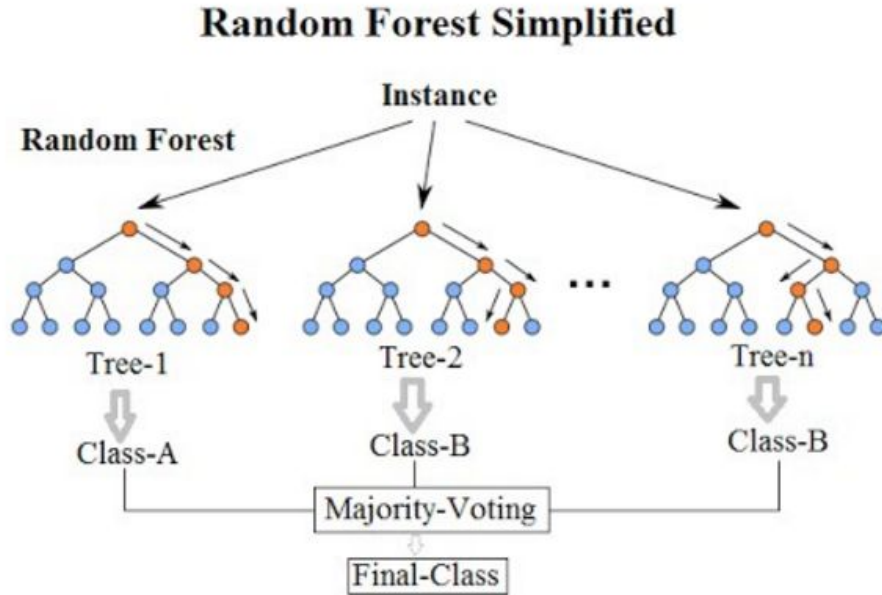
# Decision tree



# Sklearn pipeline



# Ensemble(Random forests, stacked ensemble)



\* C1, C2, and C3 are considered level 1 classifiers.

# Classification metrics(Accuracy, Precision, Recall, F1, ROC, AUC)

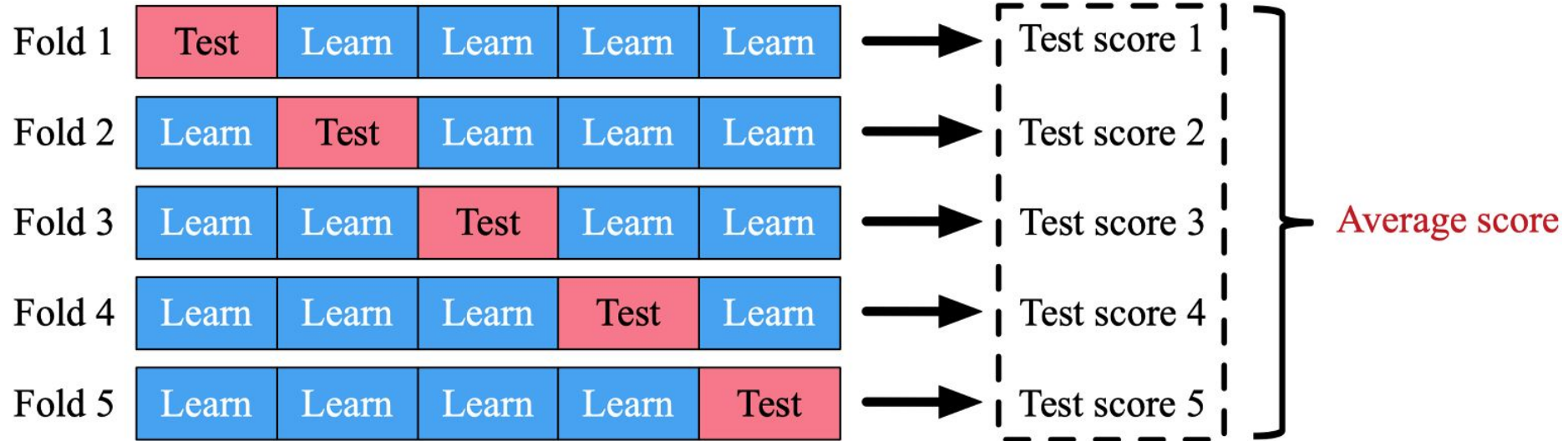
		Predicted condition		Sources: [5][6][7][8][9][10][11][12]		view · talk · edit
Total population = P + N		Predicted condition positive (PP)	Predicted condition negative (PN)	Informedness, bookmaker informedness (BM) = TPR + TNR - 1		Prevalence threshold (PT) = $\frac{\sqrt{TPR \cdot FPR} - FPR}{TPR - FPR}$
Actual condition	Actual condition positive (P)	<b>True positive (TP),</b> hit	<b>False negative (FN),</b> Type II error, miss, underestimation	True positive rate (TPR), recall, sensitivity (SEN), probability of detection, hit rate, power = $\frac{TP}{P} = 1 - FNR$		False negative rate (FNR), miss rate = $\frac{FN}{P}$ = 1 - TPR
	Actual condition negative (N)	<b>False positive (FP),</b> Type I error, false alarm, overestimation	<b>True negative (TN),</b> correct rejection	False positive rate (FPR), probability of false alarm, fall-out = $\frac{FP}{N} = 1 - TNR$		True negative rate (TNR), specificity (SPC), selectivity = $\frac{TN}{N} = 1 - FPR$
Prevalence = $\frac{P}{P + N}$		Positive predictive value (PPV), precision = $\frac{TP}{PP} = 1 - FDR$	False omission rate (FOR) = $\frac{FN}{PN} = 1 - NPV$	Positive likelihood ratio (LR+) = $\frac{TPR}{FPR}$		Negative likelihood ratio (LR-) = $\frac{FNR}{TNR}$
Accuracy (ACC) = $\frac{TP + TN}{P + N}$		False discovery rate (FDR) = $\frac{FP}{PP}$ = 1 - PPV	Negative predictive value (NPV) = $\frac{TN}{PN} = 1 - FOR$	Markedness (MK), deltaP ( $\Delta p$ ) = PPV + NPV - 1		Diagnostic odds ratio (DOR) = $\frac{LR+}{LR-}$
Balanced accuracy (BA) = $\frac{TPR + TNR}{2}$		$F_1$ score = $\frac{2 \cdot PPV \cdot TPR}{PPV + TPR} = \frac{2TP}{2TP + FP + FN}$	Fowlkes–Mallows index (FM) = $\sqrt{PPV \cdot TPR}$	Matthews correlation coefficient (MCC) = $\frac{\sqrt{TPR \cdot TNR \cdot PPV \cdot NPV} - \sqrt{FNR \cdot FPR \cdot FOR \cdot FDR}}{1}$		Threat score (TS), critical success index (CSI) = $\frac{TP}{TP + FN + FP}$

Accuracy, Precision, Recall, F1-score, ROC, AUC, ...

MAE, MSE, RMSE, R2, RMSLE, MPE, MAPE

# Cross-validation

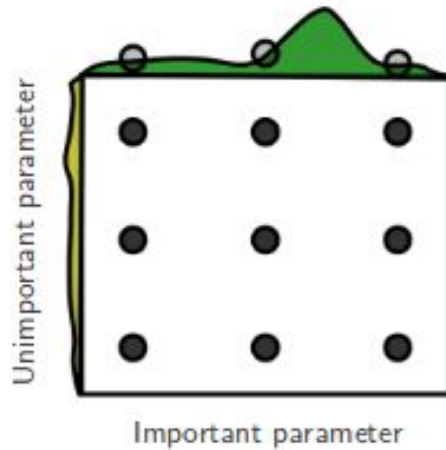
5 Cross Validation



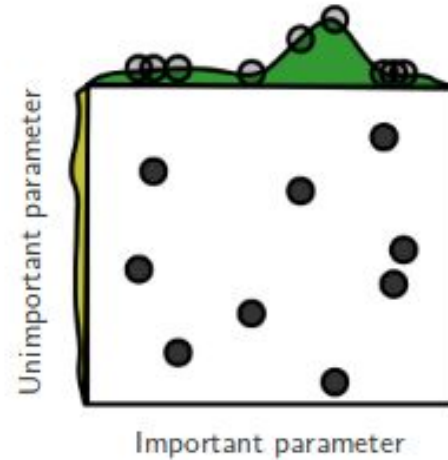


# Hyperparameter tuning

Grid Layout



Random Layout



# Sprint 3 Applied Predictive Modeling

Data science workflow

Leakage

Imbalanced data

Data wrangling

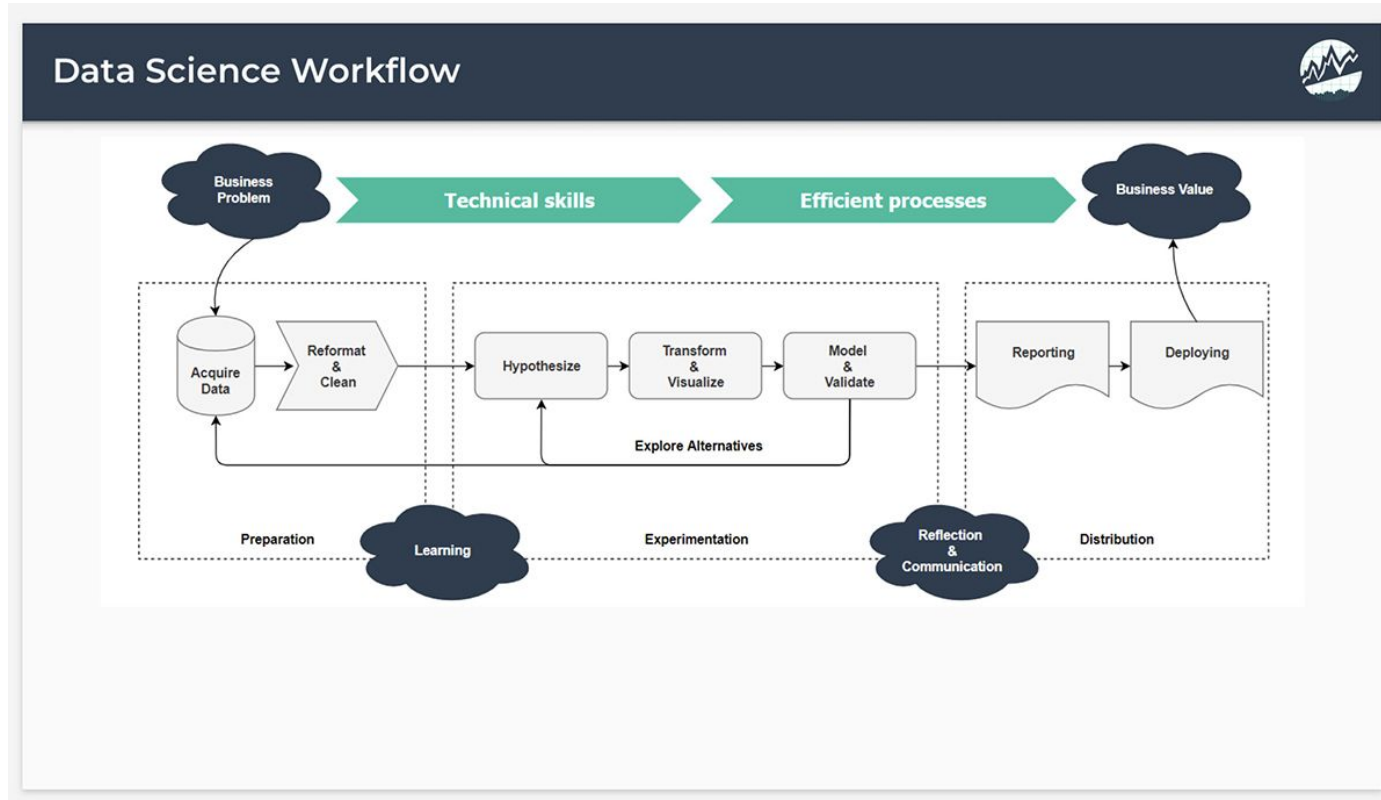
Boosting models

Feature importances

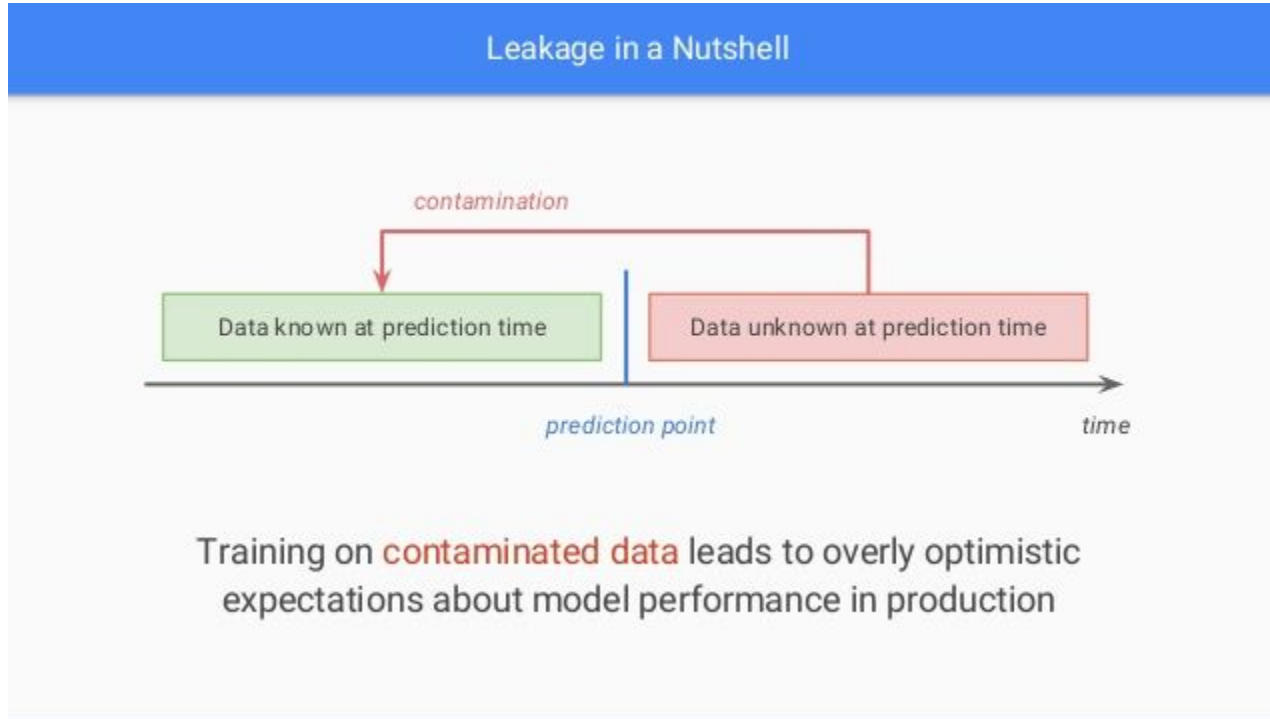
PDP

SHAP

# Data science workflow



# Leakage

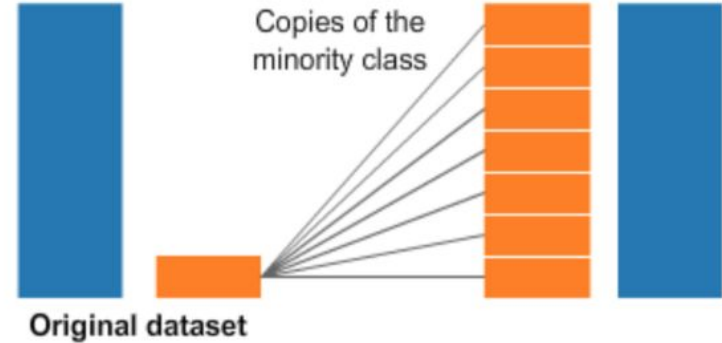


# Imbalanced data

**Undersampling**



**Oversampling**



# Data wrangling

## Data Wrangling with dplyr and tidyr Cheat Sheet



### Syntax - Helpful conventions for wrangling

#### dplyr::tbl\_df(iris)

Converts data to tbl class. tbl's are easier to examine than data frames. R displays only the data that fits onscreen.

```
Source: local data frame [158 x 5]
  Sepal.Length Sepal.Width Petal.Length
1           5.1           3.5           1.4
2           4.9           3.0           1.4
3           4.7           3.2           1.3
4           4.6           3.1           1.5
5           5.0           3.6           1.4
...
Variables not shown: Petal.Width (dbl),
Species (fctr)
```

#### dplyr::glimpse(iris)

Information dense summary of tbl data.

#### utils::View(iris)

View data set in spreadsheet-like display (note capital V).

#### dplyr::%>%

Passes object on left hand side as first argument (or argument) of function on righthand side.

$x \%>\% f(y)$  is the same as  $f(x, y)$   
 $y \%>\% f(x, \dots, z)$  is the same as  $f(x, y, z)$

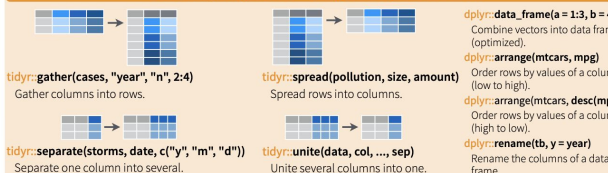
"Piping" with %>% makes code more readable, e.g.

```
iris %>%
  group_by(Species) %>%
  summarise(avg = mean(Sepal.Width)) %>%
  arrange(avg)
```

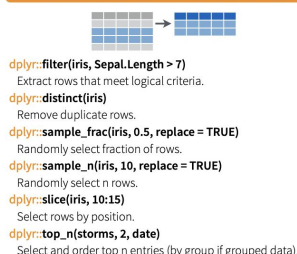
### Tidy Data - A foundation for wrangling in R



### Reshaping Data - Change the layout of a data set

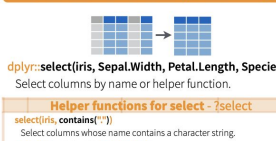


### Subset Observations (Rows)



Logic in R - ?Comparison, ?base::Logic			
<	Less than	!=	Not equal to
>	Greater than	%in%	Group membership
==	Equal to	is.na	Is NA
<=	Less than or equal to	is.na	Is not NA
>=	Greater than or equal to	&,  , !, xor, any, all	Boolean operators

### Subset Variables (Columns)



Helper functions for select - ?select

select(iris, contains(""))  
Select columns whose name contains a character string.

select(iris, ends\_with("Length"))  
Select columns whose name ends with a character string.

select(iris, everything())  
Select every column.

select(iris, matches("L"))  
Select columns whose name matches a regular expression.

select(iris, num\_range("x", 1:3))  
Select columns whose name matches a regular expression.

select(iris, x1, x2, x3, x4, x5)  
Select columns named x1, x2, x3, x4, x5.

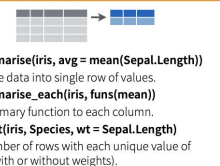
select(iris, one\_of("Species", "Genus"))  
Select columns whose names are in a group of names.

select(iris, starts\_with("Sepal"))  
Select columns whose name starts with a character string.

select(iris, Sepal.Length:Petal.Width)  
Select all columns between Sepal.Length and Petal.Width (inclusive).

select(iris, -Species)  
Select all columns except Species.

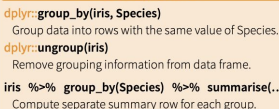
### Summarise Data



Summarise uses **summary functions**, functions that take a vector of values and return a single value, such as:

**first** First value of a vector.  
**last** Last value of a vector.  
**nth** Nth value of a vector.  
**n** # of values in a vector.  
**n\_distinct** # of distinct values in a vector.  
**IQR** IQR of a vector.  
**min** Minimum value in a vector.  
**max** Maximum value in a vector.  
**mean** Mean value of a vector.  
**median** Median value of a vector.  
**var** Variance of a vector.  
**sd** Standard deviation of a vector.

### Group Data



iris %>% group\_by(Species) %>% summarise(...)  
Compute separate summary row for each group.



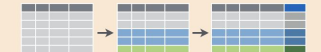
### Make New Variables



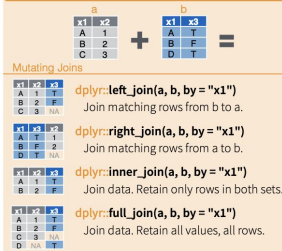
**lead** Copy with values shifted by 1.  
**lag** Copy with values lagged by 1.  
**dense\_rank** Ranks with no gaps.  
**min\_rank** Ranks. Ties get min rank.  
**percent\_rank** Ranks rescaled to [0, 1].  
**row\_number** Ranks. Ties get to first value.  
**ntile** Bin vector into n buckets.  
**between** Are values between a and b?  
**cume\_dist** Cumulative distribution.

**cumall** Cumulative all  
**cumany** Cumulative any  
**cummean** Cumulative mean  
**cumsum** Cumulative sum  
**cummax** Cumulative max  
**cumin** Cumulative min  
**prod** Element-wise prod  
**max** Element-wise max  
**min** Element-wise min

iris %>% group\_by(Species) %>% mutate(...)  
Compute new variables by group.



### Combine Data Sets



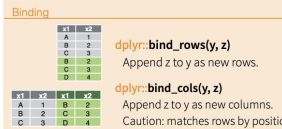
**Filtering Joins**

**semi\_join(a, b, by = 'x1')**  
All rows in a that have a match in b.  
**anti\_join(a, b, by = 'x1')**  
All rows in a that do not have a match in b.

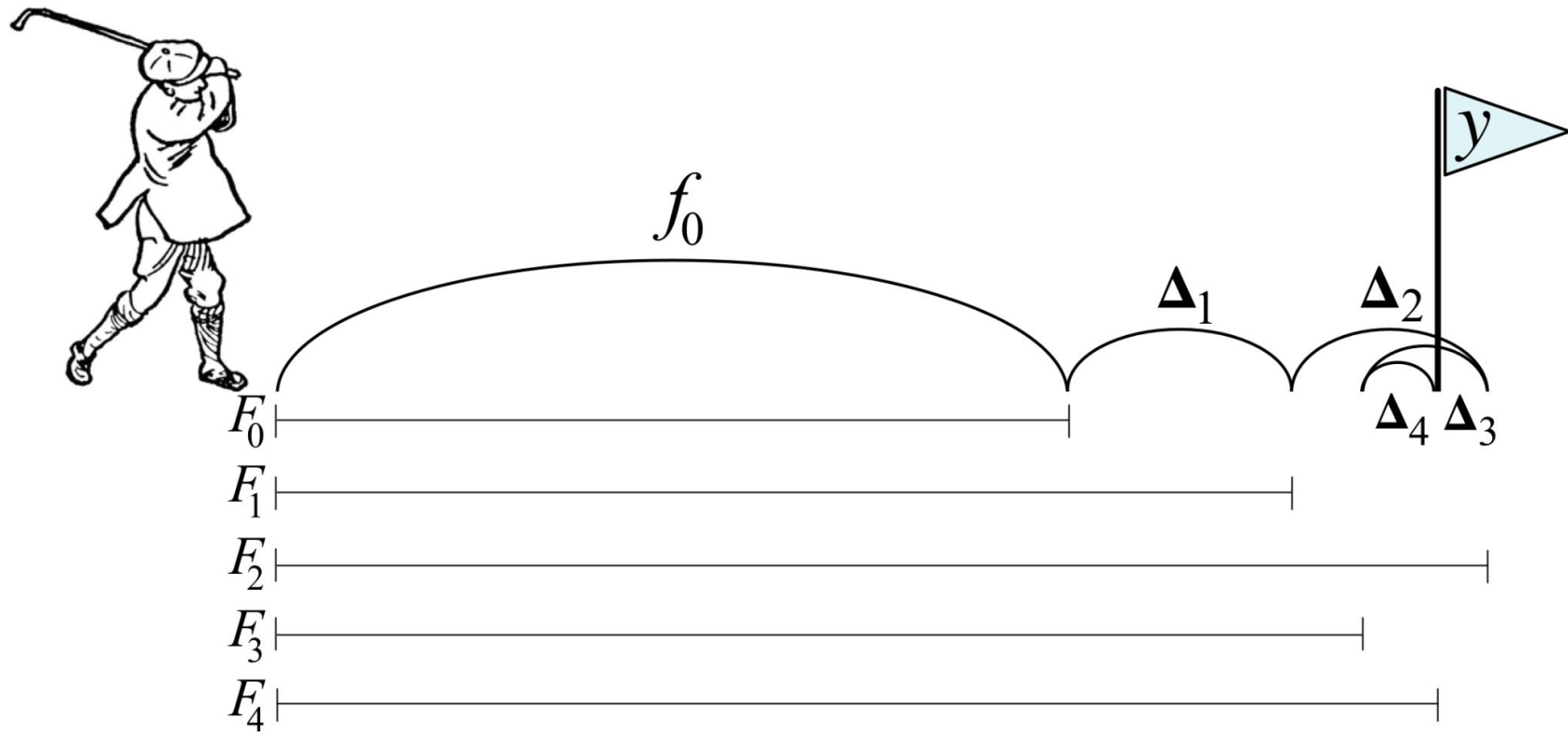


**Set Operations**

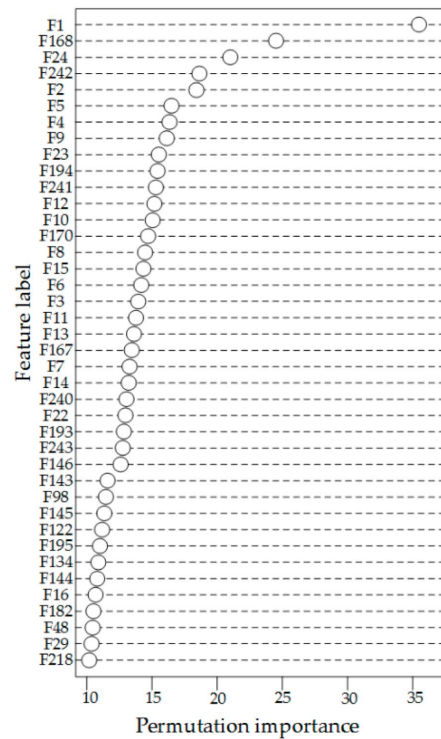
**intersect(y, z)**  
Rows that appear in both y and z.  
**union(y, z)**  
Rows that appear in either or both y and z.  
**setdiff(y, z)**  
Rows that appear in y but not z.



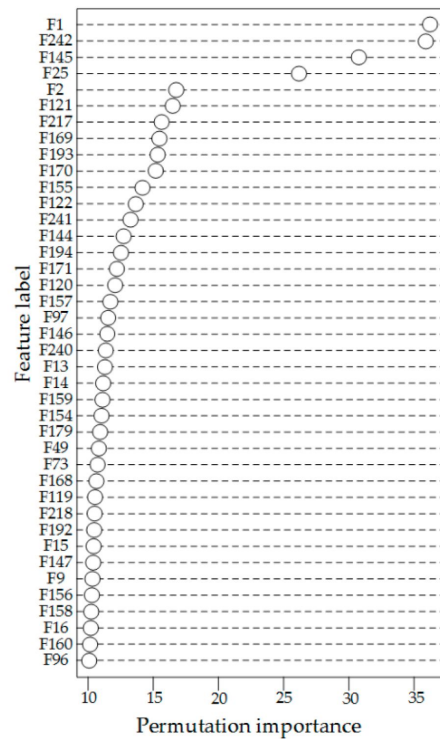
# Boosting model



# Feature importances



(a)

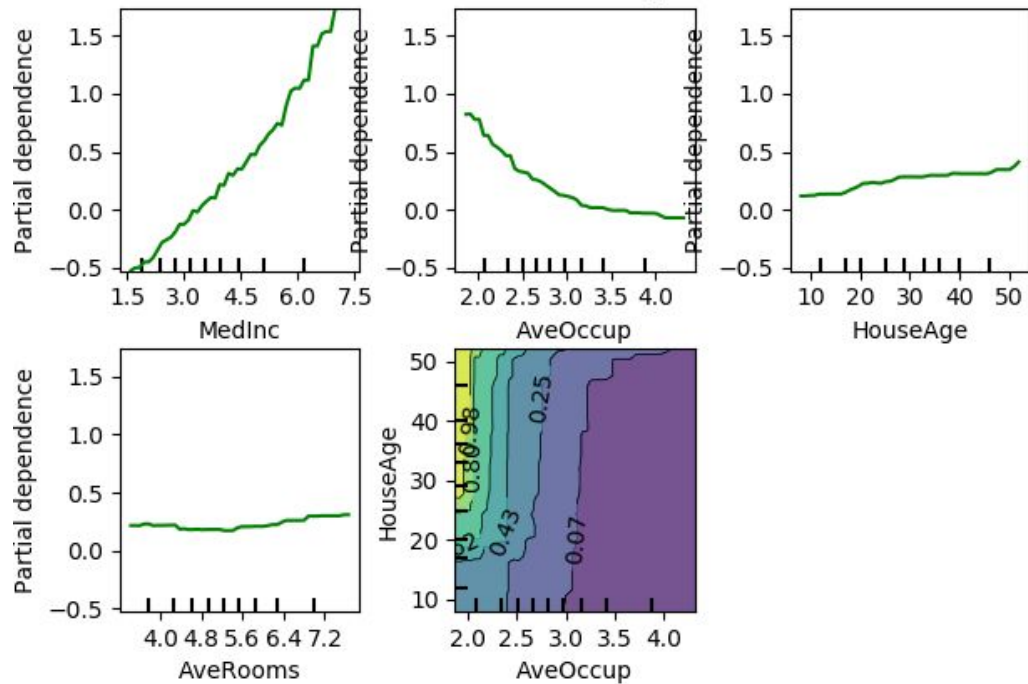


(b)

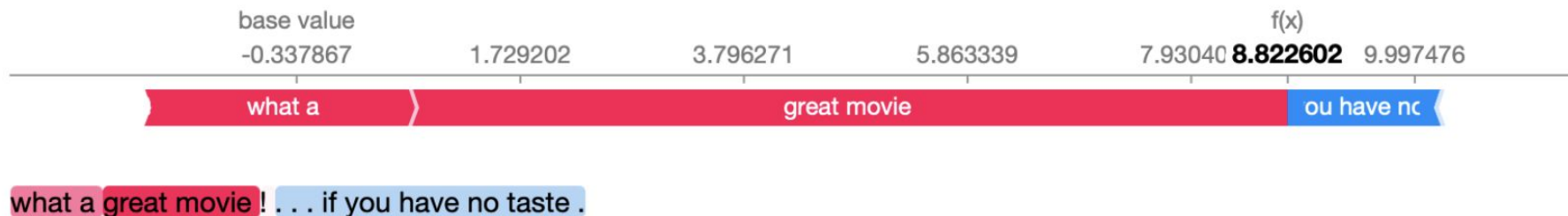
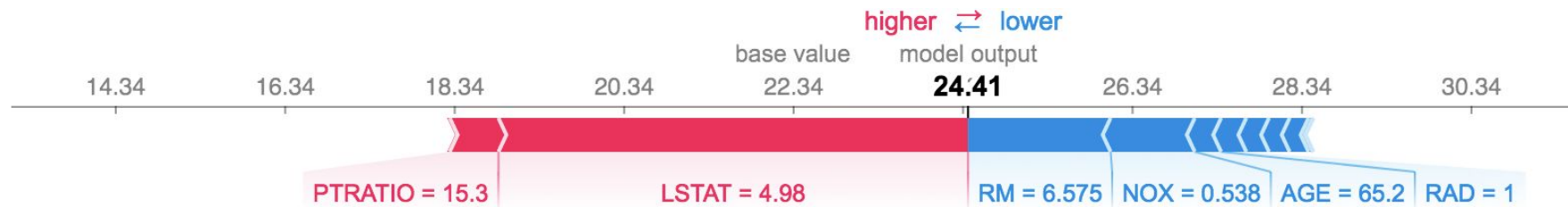


# PDP

Partial dependence of house value on nonlocation features  
for the California housing dataset



# SHAP



# Deep Learning

Image recognition

Natural language processing

Speech recognition

Transfer learning

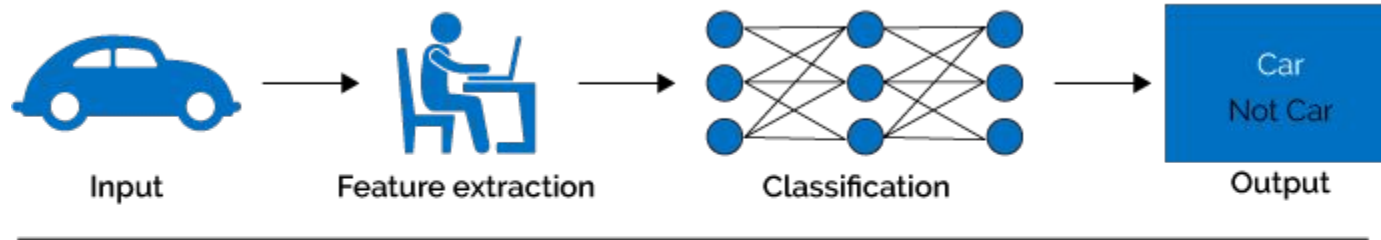
Semi-Supervised learning

Unsupervised learning, Generative model(GAN, VAE)

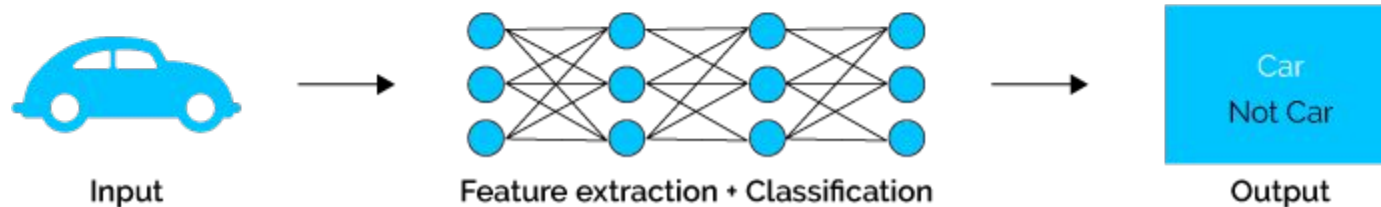
....

# Machine learning vs Deep learning

## Machine Learning



## Deep Learning



# Machine learning algorithms

