



Warasinee Chaisangmongkon, PhD

Lecturer at Institute of Field Robotics, KMITT; Data Scientist at Big Data Experience Center

- PhD from Yale University,
- 10 years experiences in machine learning & data mining
- Lead scientist of corporate data science projects (banking, business digitizations)
- Founded IDEA LAB : R&D in machine learning for businesses



Introduction to Machine Learning

.....

Instructor: Warasinee Chaisangmongkon, PhD

Day 2 : Predictive Modeling

2

- ④ What is Machine Learning
- ④ Feature, Target, Sample
- ④ Model & Cost Function
- ④ Overfitting

INTRODUCTION TO MACHINE LEARNING

3

Introduction to Machine Learning

What it means to make predictive models for your data.

Decision Trees and Random Forest Models

Interactive lessons in tree-based modeling.

Feature Engineering

How to get your data ready for modeling and how to select the best features.

Lab

Use your knowledge to fit a model to predict house prices!

Linear Regression and Logistic Regression

Interactive lessons in linear and logistic regression in R

Conclusions

Other regression and classification models.

K-Nearest Neighbor

Interactive lessons on KNN methods.

4

THE MOST BASIC UNDERSTANDING



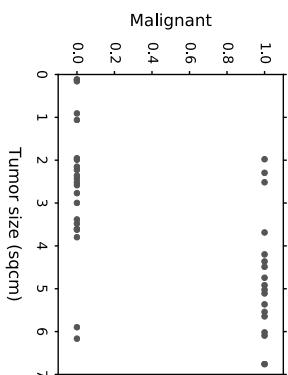
- It's all about letting computer learns what 'input' is associated to what 'output'.

- Example: given behavioral profiles of a customer A, algorithms predict the chance this customer is going to leave.
- Example: given inputs from sensors and cameras, the robotic algorithm pushes out the appropriate movement.

INTELLIGENT SYSTEM WITH MACHINE LEARNING

- People provide Facebook the images and tags of names in the photos.
- Over time Facebook learned to associate names with faces and can automatically recognize these people.

A SIMPLE EXAMPLE



- A doctor has seen hundreds of patients with tumors. They measured the size of tumors and test whether they are malignant.
- The historical data let the doctor make prediction that bigger tumor is more likely to be malignant.

THE MOST BASIC EXAMPLE

WHAT IS A MODEL

i	x	y
1	Size (m ²) (Mbaht)	Price (Mbaht)
2	50 = x_1	1.4 = y_1
3	128 = x_2	2.6 = y_2
4	24 = x_3	0.8 = y_3
i	$\dots = x_i$	$\dots = y_i$

In general,

x: feature (input)

y: target (output)

i: sample index

After we have a model, we can predict y value from any x value

$$\text{Model: } h(x_i) = 4.717 + 0.038 * (x_i)$$

Notice that:
 $h(x)$ and y are not the same
 $h(x)$: value we predict
 y : the actual value

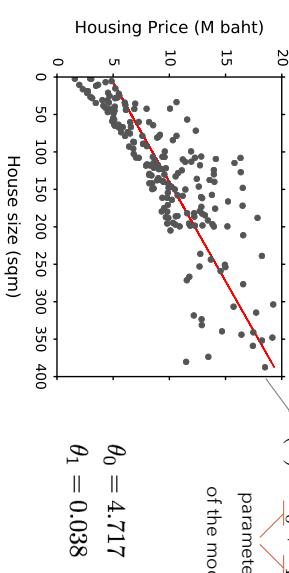


THE MOST BASIC EXAMPLE

WHAT IS A MODEL

i	x	y
1	Size (m ²) (Mbaht)	Price (Mbaht)
2	50	1.4
3	128	2.6
4	24	0.8
i	\dots	\dots

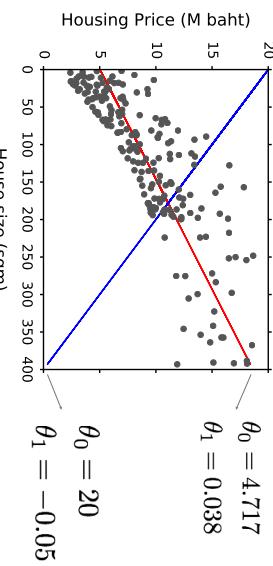
An agent has been selling 300 houses in the last years and want to be able to predict the price of a house by just knowing the size of the house.



A model is a function that takes an input and yields the values we want to predict.

GOOD AND BAD MODELS

Model: $h(x) = \theta_0 + \theta_1x$



WHAT IS A MODEL

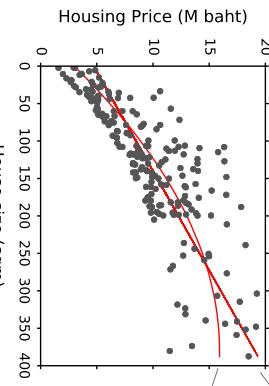
For different models you have different functions, with perhaps different number of parameters. Different models can be fitted to the same data set.

$$h(x) = \theta_0 + \theta_1x$$

$$h(x) = \theta_0 + \theta_1x + \theta_2x^2$$

Model is a function of both features (x) and parameter (theta).

$$h(x; \theta)$$



LOWERING COST FUNCTION

Model: $h(x) = \theta_0 + \theta_1x$

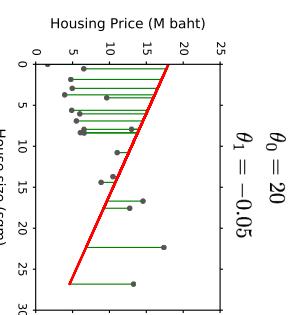
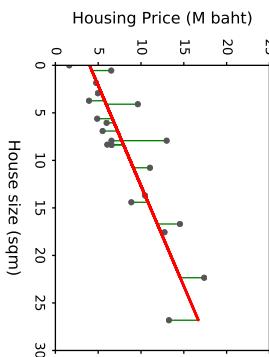


COST FUNCTION

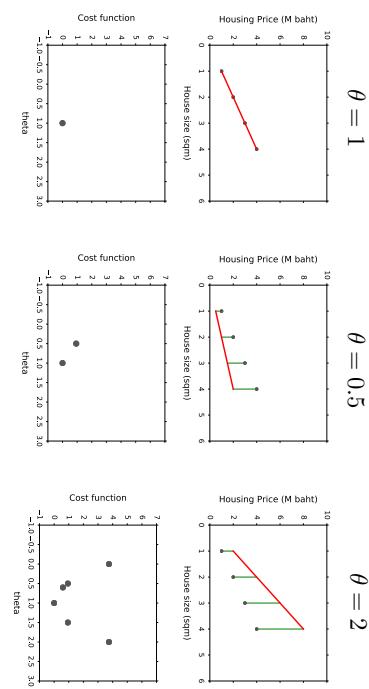
Cost function (or loss function) is a measure of whether a model is a good fit to the data.

For example, a famous cost function is called 'squared error' function that takes the difference between what you predict and the actual data value and square it.

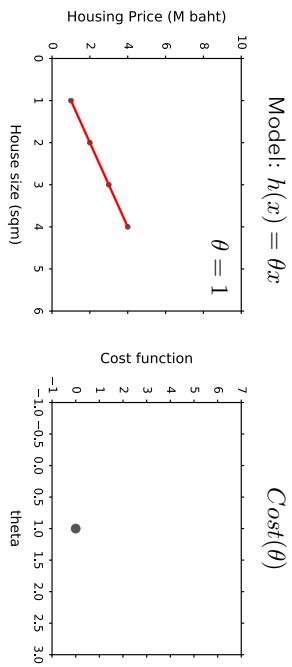
$$\sum_i (h(x_i) - y_i)^2$$



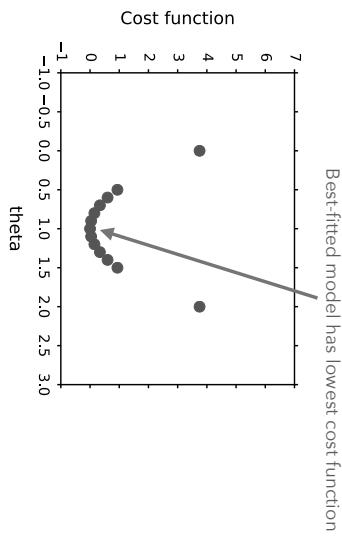
MINIMIZING COST FUNCTION



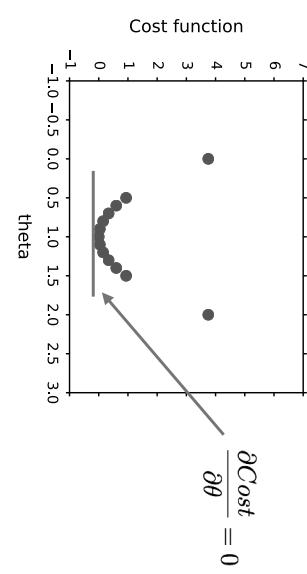
MINIMIZING COST FUNCTION



MINIMIZING COST FUNCTION



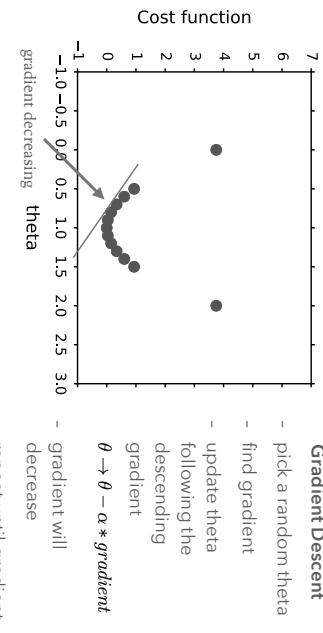
Best-fitted model has lowest cost function



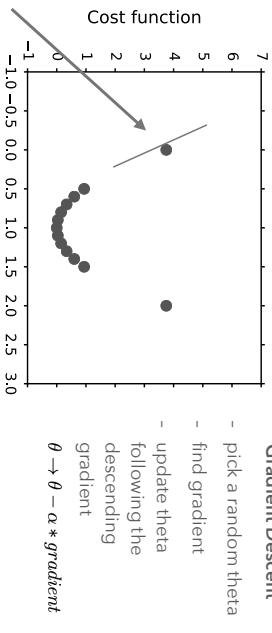
Some optimization algorithm is so simple,
you just solve an equation and done.

Training machine learning models = minimizing cost function
Often accomplished by using 'optimization algorithms'

MINIMIZING COST FUNCTION



MINIMIZING COST FUNCTION



TAKE-HOME MESSAGES

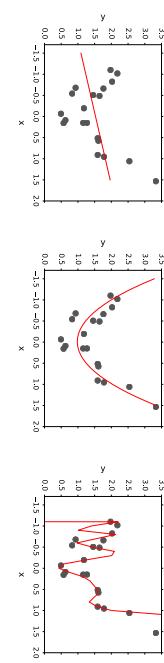


Gradient Descent

- pick a random theta
- find gradient
- update theta following the descending gradient
- $\theta \rightarrow \theta - \alpha * gradient$

- Machine learning is a science of letting computers take input (features), pass them through a model (an equation or a system of equations) and return a predicted output.
- Models consist of parameters we can adjust to fit the data.
- The goodness of fit is judged by a cost function which measures how much our predictions are far away from actual data (target).
- Machine learning application always involves an algorithm which optimizes cost function.

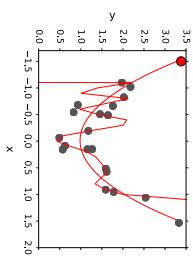
OVERFITTING



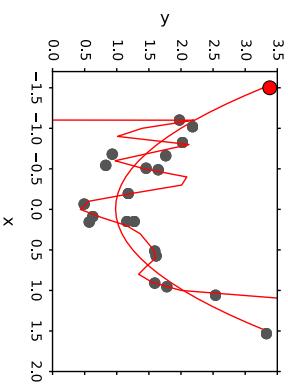
FACTS ABOUT OVERFITTING



- If the number of features (n) is really large, you can fit y with really high precision.
- The more data you have (compared to parameters) the less likely your model will overfit, because noise will be more likely to average out.



OVERFITTED MODELS CANNOT GENERALIZE

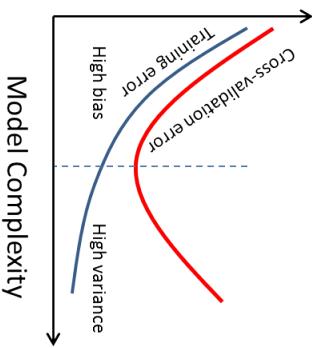


FACTS ABOUT OVERFITTING



- The more parameters, the larger solution space you fit to the data, and this is not a good thing.
- The less parameters, the more likely your model will underfit (not having enough feature to make predictions about the target).
- This is called bias-variance tradeoff (bias: model has bias, i.e. too strong assumption) (variance: too much errors/noise is fitted).

Mean Error



FACTS ABOUT OVERFITTING

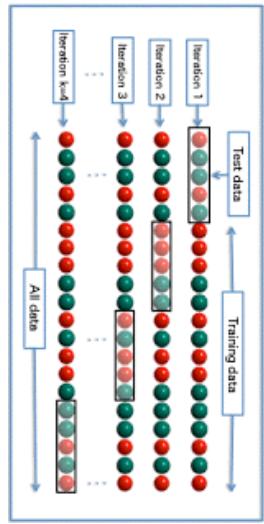


- If the number of features (n) is really large, you can fit y with really high precision.
- The more data you have (compared to parameters) the less likely your model will overfit, because noise will be more likely to average out.

AVOID OVERFITTING



- ▶ Cross-validate your models: train models with one set of data (training set) and test them with another set of data (test set)



AVOID OVERFITTING



- ▶ Cross-validate your models: train models with one set of data (training set) and test them with another set of data (test set)

cross validation

70% of data were used for training the algorithm
30% of data preserved for testing

REGULARIZATION



- ▶ Regularization is a mathematical way to reduce overfitting automatically, by limiting the influence of each feature.

$$\sum_i (h(x_i) - y_i)^2 + \lambda \sum_j \theta_j^2$$

The normal cost function
norm of parameter vector

- ▶ To minimize this cost function you have to minimize both first and second terms
- ▶ Minimize the second term means less overfitting.

AVOID OVERFITTING



- ▶ Reduce the number of features and parameters

- ▶ Manual selection
- ▶ Algorithmic selection (information gain, pruning, stepwise algorithm)
- ▶ Regularization

AVOID OVERFITTING

36

- ▶ Cross-validate your models: train models with one set of data (training set) and test them with another set of data (test set)
- ▶ Reduce the number of features and parameters
- ▶ Manual selection
- ▶ Algorithmic selection (information gain, pruning, stepwise algorithm)
- ▶ Regularization: cost function is split to two parts

$$\sum_i (h(x_i) - y_i)^2 + \lambda \sum_j \theta_j^2$$

REGULARIZATION

$$\sum_i (h(x_i) - y_i)^2 + \lambda \sum_j \theta_j^2$$

The normal cost function norm of parameter vector

FEATURE ENGINEERING

- ④ What is it and why
- ④ Data preprocessing
- ④ Preprocessing w/t R

- ▶ Lambda is called ‘regularization parameter’
- ▶ Because we adjust lambda to adjust the weight of the two cost function terms
- ▶ High lambda: severely limit parameter size
- ▶ Low lambda: allow parameters to scale up more freely, give more importance to lowering error

Feature engineering is the process of transforming raw data into features that better represent the underlying problem to the predictive models, resulting in improved model accuracy.

35

FEATURE ENGINEERING



- ▶ Preprocessing (filtering, cleaning, linking)
- ▶ Feature extraction
- ▶ Create a new feature from existing feature
 - ▶ Create a new feature from 2 or more features
- ▶ Feature selection
- ▶ Dimensionality reduction
- ▶ Deep learning (a.k.a feature learning)

FEATURE SCALING



- ▶ Imagine if x_1 and x_2 are not similar in scale
- ▶ For example
 - ▶ x_1 = number of bedrooms (0-20)
 - ▶ x_2 = area of the house (24-3000 sqm)
- ▶ This means the scale of your theta will also be different.

DATA SCIENTIST TASKS



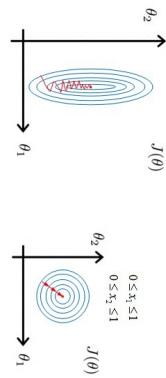
- ▶ 50% - Obtaining data and cleaning data
- ▶ 30% - Feature engineering
- ▶ 20% - Modeling

COMMON PREPROCESSING OPERATIONS



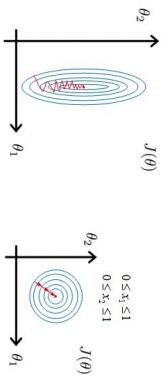
- ▶ Feature scaling
- ▶ Dealing with skewed data
- ▶ Dealing with missing data

FEATURE SCALING



- ▶ Simple solution would be to scale the features.
 - ▶ x_1 : number of bedroom $\rightarrow x_1 = \#bedroom / \max \#bedroom$
 - ▶ x_2 : area of the house $\rightarrow x_2 = \text{area} / \max \text{area}$
 - ▶ Try to get features to stay within -1 to 1

FEATURE SCALING



- ▶ A symmetric contour
 - ▶ Landscape of cost function is like ગુમારાં
 - ▶ Gradient descent algorithm has a hard time with such surface

FEATURE SCALING

- These three techniques are suitable for most problems

normalization with max, min, mean
standardization
a.k.a. z-score

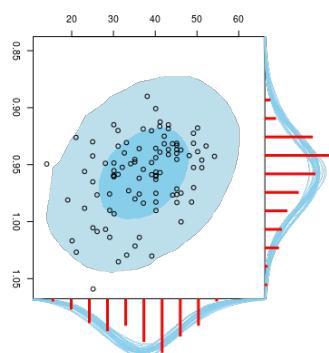
$$x_1 := \frac{x_1}{\max(x_1)}$$

$$x_1 := \frac{x_1 - mean(x_1)}{max(x_1)}$$

$$x_1 := \frac{x_1 - \text{mean}(x_1)}{\text{std}(x_1)}$$

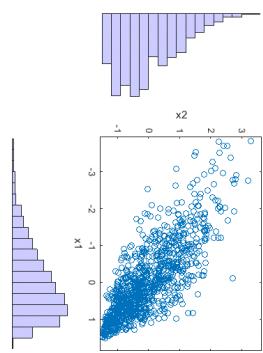
What's the range of these rescaled variables?

SKEWED DATA



Normally, we operate regression with the assumption that variables are normally distributed.

SKewed DATA



Regression does not guarantee solution for skewed distribution.

SKewed DATA

- ▶ Regression has the following assumptions:

- ▶ Linear relationship
- ▶ Multivariate normality

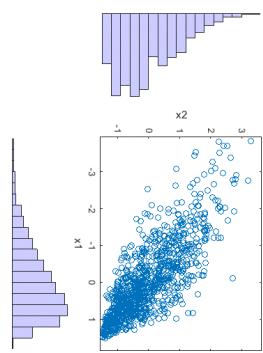
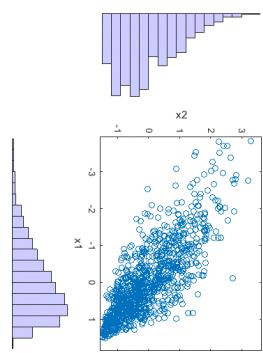
- ▶ No or little multicollinearity

- ▶ No auto-correlation

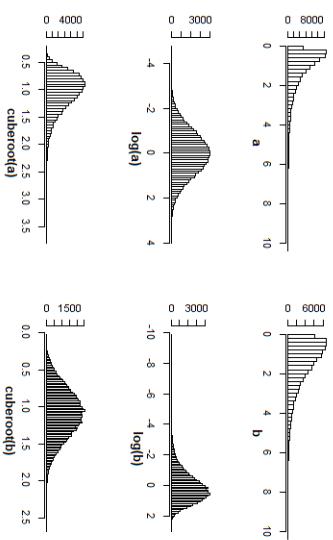
- ▶ Homoscedasticity (all variables have same variance)

MISSING DATA

- ▶ If data are assumed to be missing at random, we may simply ignore the data
- ▶ You went to all houses and randomly for some houses, you took time to do detailed house measurement
- ▶ In this case, you simply remove the missing data from the analysis (cut the whole row)
- ▶ Be aware that missing data might not be as random as you think



SKewed DATA



MISSING DATA

- ▶ Mean or mode substitution
- ▶ Missing income? Fill it with average income?
- ▶ Don't know if the patient is left-handed or right-handed, assume right-handed, because it's more common.
- ▶ Weaken correlation and covariance.

i	Size (m ²)	Price (Mbaht)	# Bed	Price
1	50	1.4	2	1.4
2	128	2.6	?	2.6
3	24	0.8	1	0.8
4	?	1.2	2	1.2
i

MISSING DATA

i	Size (m ²)	Price (Mbaht)	# Bed	Price
1	50	1.4	2	1.4
2	128	2.6	?	2.6
3	24	0.8	1	0.8
4	?	1.2	2	1.2
i

MISSING DATA

- ▶ Dummy variable
- ▶ Customers are divided to high (3), medium (2), low income (1), assume that customers with missing income is of category 0.
- ▶ In some cases, this is perfect because people that avoid filling in income might have interesting characteristics.
- ▶ In some cases, this method is not so good because you introduced and extra value that is not driven by fact.
- ▶ Listwise means cutting the whole row (reduced sample)
- ▶ Pairwise means cutting only the missing value (can't do multivariate analysis)
- ▶ For example, if customers forgot to fill in their ages, you would not assume all the missing data have something in common.

Big Data

- Most problems you will face in the real world is gigantic.
- Millions of rows
- Hundreds or thousands of features
- Your algorithm will take forever to run
- What can we do about it?
- We might be able to look through all the features and manually select them.
- But that would waste so much time and resources
- So maybe do automated feature selection?

What Feature Matters Most?

- The algorithm predicts whether the email is spam or not. Which feature is most useful for the prediction?
- Feature 1: whether the email contains the word 'viagra'
- Feature 2: whether the email is sent from a Nigerian Prince
- Feature 3: whether the email is sent from one person to a massive amount of people
- For all 1000 features you have calculated, maybe only a few features are important.
- Feature selection algorithm gives you insight and interpretability of your model.

53

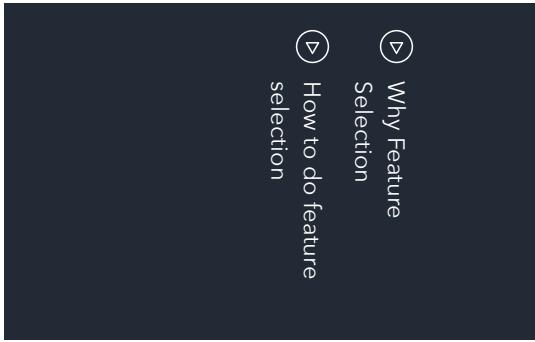
Feature Selection:
**The process of selecting
the most relevant features
to be included in**

- ④ Why Feature Selection
- ④ How to do feature selection

FEATURE SELECTION

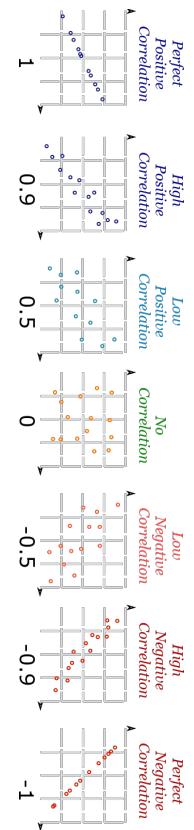
55

the machine learning model



Curse of Dimensionality

- This is one of the most important problems in machine learning.
- Take linear regression for example
- The more features you have, the more parameters you need to fit the model
- If you have 2 features, your solution space has 2 dimensions (small possible values)
- If you have 1000 features, your solution space has 1000 dimensions (huge amount of possible values).
- Your algorithm can take a lot of time to find solution.



How to find a good feature

- You need to find features that have high correlation to your target.
 - Don't care whether it's a positive or negative correlation
 - The larger the number the better

Curse of Dimensionality

- This is one of the most important problems in machine learning.

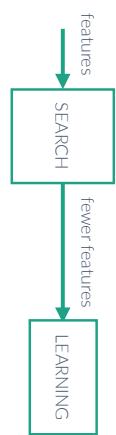
- Imagine you have a large amount of features, each can have infinite number of values.

- You will need an enormous amount of training data is required to ensure that there are several samples with each combination of values.

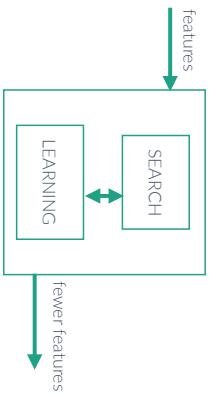
- If you have limited samples, which do not cover the whole space, your model loses predictive power.

Feature Selection

Filtering Methods



Wrapping Methods



Analysis of Variance

- Classification problem: Y can only be class 0 or class 1. Find variance of X within class and between classes.

$$F\text{-Value} = \frac{\text{Variance between classes}}{\text{Variance within class}}$$

V between class is **high**
V within class is **low**
F-Value is **high**
Feature X is **important**

V between class is **low**
V within class is **high**
F-Value is **low**
Feature X is **not important**

How to find a good feature

- Correlation is not the only measure that tells you 'how much x is related to y'
there are other measures we use.
- Such as:
- ANOVA (Analysis of Variance)
- Chi2
- Mutual Information

Weight	Class	Weight	Class	Weight	Class
50	Adult	68	Thailand	65	Human
80	Adult	75	China	1000	Animal
12	Children	80	Thailand	0.1	Animal
30	Children	82	China	60	Animal
...	30	Human

Analysis of Variance Quiz

V between class is ...
V within class is ...
F-Value is ...
Feature is ...

V between class is ...
V within class is ...
F-Value is ...
Feature is ...

V between class is ...
V within class is ...
F-Value is ...
Feature is ...