

Data Analysis and Interpretation

Statistical Data Analysis

Recap from Week 3

- ❖ Data Distribution
- ❖ Data Relationship
- ❖ Sequential Data
- ❖ Spatial Data
- ❖ Temporal Data
- ❖ Bar chart
- ❖ Pie chart
- ❖ Scatter plot
- ❖ Heat map
- ❖ Box-and-whisker plot
- ❖ Transition graph
- ❖ Map
- ❖ Time Series

3803ICT course structure

W1. Introduction to Data Analytics

Data Preparation and Preprocessing

W2. Data Preparation and Preprocessing

Data Analysis and Interpretation

W3. Exploratory Data Analysis

W4. Statistical Data Analysis

W5. Predictive Data Analysis

Visualization

W6. Visual Data analysis

Analysis of special types of data

W7. Time Series

W8. Textual Data

W9. Graph Data

Analysis with big data infrastructure

W10. Distributed Data Analysis

W11. Cloud-based Data Analysis

W12. Revision

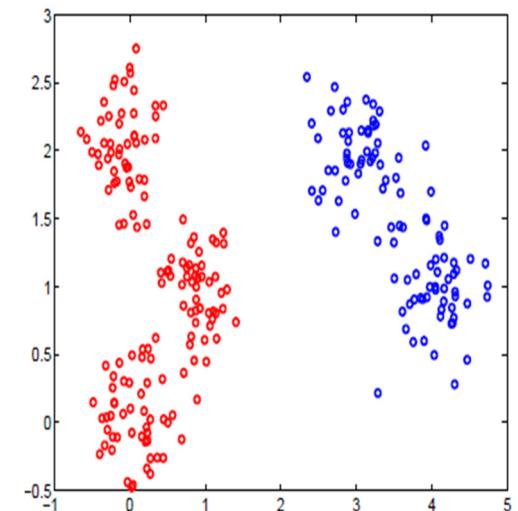
Data Modeling - Statistical Analysis

Revise the content of 2802ICT Intelligent Systems and focus on data modeling

- I. Clustering – Unsupervised Learning
- II. Classification – Supervised Learning
- III. Regression – Supervised Learning
- IV. Model Evaluation

I. Clustering

- ❖ Cluster: A collection/group of data objects/points
 - similar (or related) to one another within the same group
 - dissimilar (or unrelated) to the objects in other groups
- ❖ Cluster analysis
 - find *similarities* between data according to characteristics underlying the data and grouping similar data objects into clusters
 - unsupervised learning: no predefined classes for a training data set
 - Two general tasks: identify the “natural” clustering number and properly group objects into “sensible” clusters
- ❖ Why clustering?
 - A stand-alone tool to gain an insight into data distribution
 - A pre-processing step of other data analysis routines



how many clusters?

Clustering: Applications

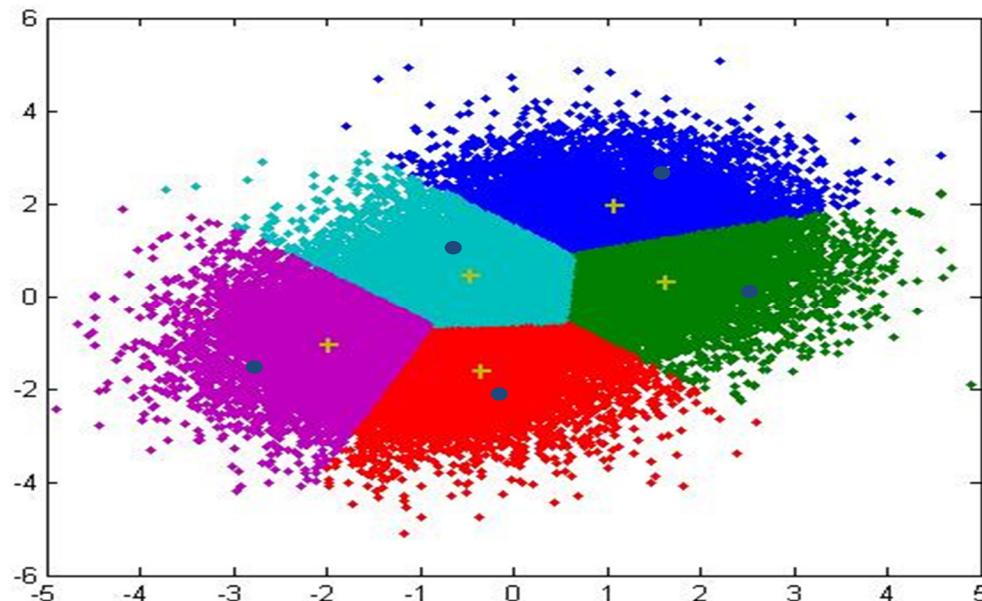
- ❖ Applications:
 - **Bank Security:** fraud pattern discovery
 - **Biology:** taxonomy of living things such as kingdom, phylum, class, order, family, genus and species
 - **City-planning:** Identifying groups of houses according to their house type, value, and geographical location
 - **Climate change:** understanding earth climate, find patterns of atmospheric and ocean
 - **Finance:** stock clustering analysis to discover correlation between shares
 - **Information retrieval/organisation:** Google search, topic-based news
 - **Land use:** Identification of areas of similar land use in an earth observation database
 - **Marketing:** Help marketers discover distinct groups in their customer bases, and then use this knowledge to develop targeted marketing programs
 - **Social network mining:** special interest group automatic discovery (**community recognition**)

Clustering: Methods

- ❖ Partitioning approach (1 variable)
 - Construct various partitions and then evaluate them by some criterion, e.g., minimizing the sum of square distance
 - Typical methods: **K-means**, K-medoids, CLARANS,
- ❖ Hierarchical approach (1 variable)
 - Create a hierarchical decomposition of the set of data (or objects) using some criterion
 - Typical methods: **Agglomerative**, Diana, Agnes, BIRCH, ROCK,
- ❖ Density-based approach (1 variable)
 - Based on **connectivity** and **density** functions
 - Typical methods: DBSCAN, OPTICS, DenClue,
- ❖ Co-clustering (2 variables)

K-means clustering

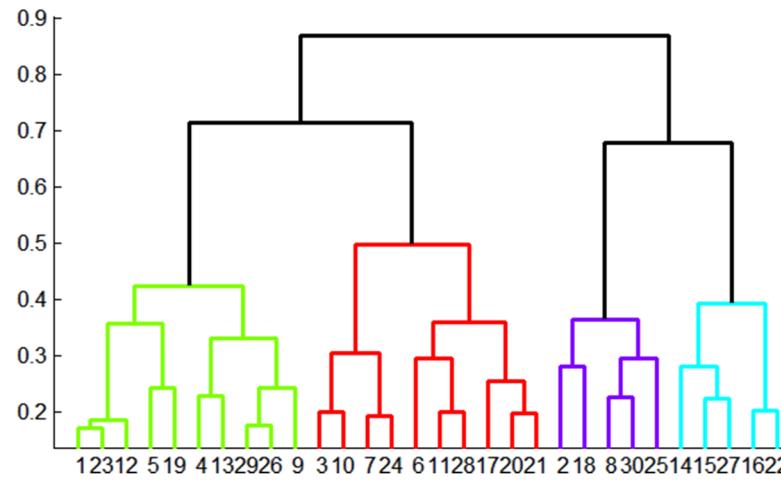
- ❖ A simple greedy algorithm (usually called Lloyd's algorithm)
- ❖ Locally **minimizes the distance** (usually Euclidean) from data points to their respective centroids:
 - Find the closest cluster centroid for each item, and assign it to that cluster.
 - Recompute the cluster centroid as the mean of items, for the newly assigned items in the cluster.



An iteration of k-means clustering (old centers = cross, new centers = dot)

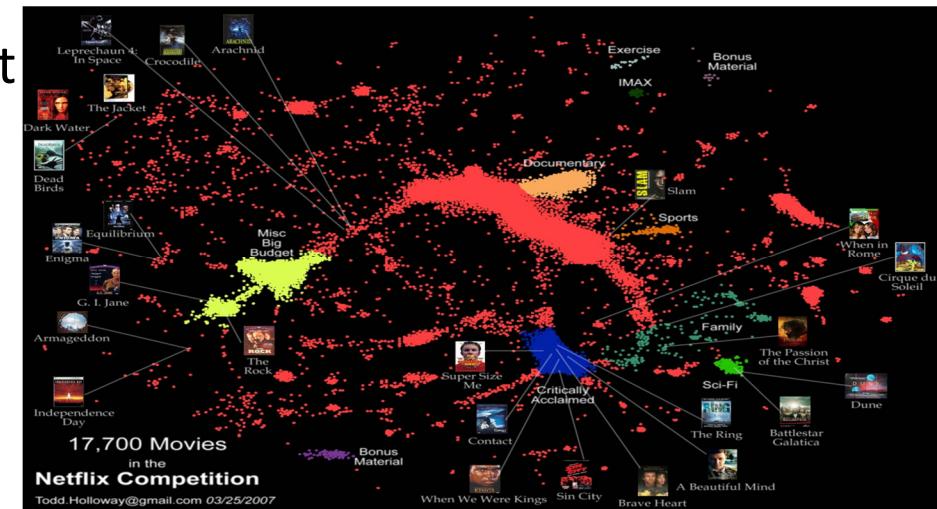
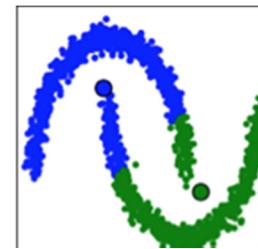
Hierarchical Clustering

- ❖ Agglomerative (bottom up):
 - Initially, each point is a cluster
 - Repeatedly combine the two “nearest” clusters into one
- ❖ Divisive (top down):
 - Start with one cluster and recursively split it
- ❖ Does not require the number of clusters k in advance
- ❖ But needs a cut-off threshold to return clusters



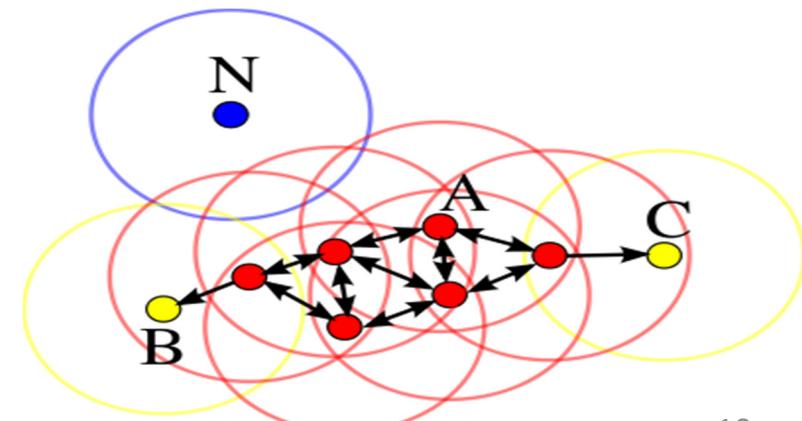
Density-based clustering: DBScan

- ❖ **Motivation:** Centroid-based clustering methods like k-means favor clusters that are spherical, and have great difficulty with anything else, e.g.



- ❖ **How it works:** follows the shape of **dense neighborhoods** of points.

- Core points have at least minPts neighbors in a sphere of diameter ϵ around them.
- The red points here are core points with at least minPts = 3 neighbors in an ϵ -sphere around them.

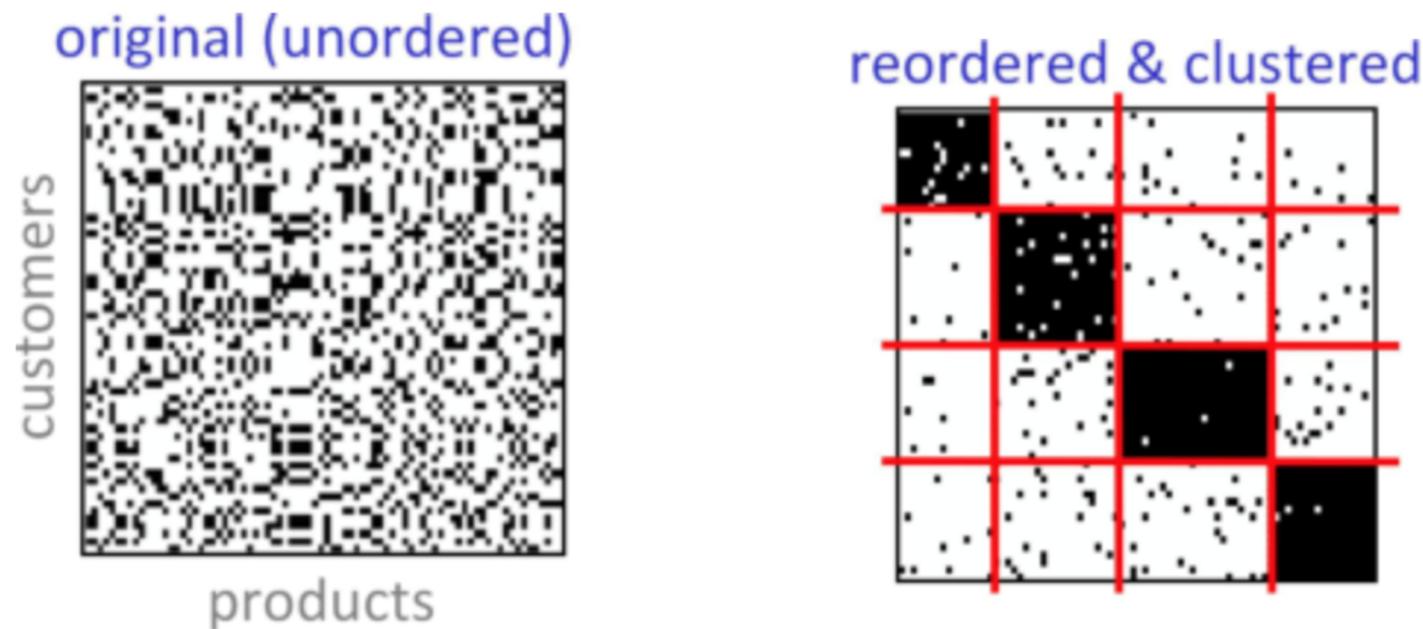


Co-clustering

- ❖ Real world data is often bimodal
 - Created by a **joint interaction** between two variables
 - E.g. a customer buying a product
- ❖ Traditional clustering (1 variable):
 - Customer clustering: group customers buying the same products, **regardless of** product clusters
 - Product clustering: group products being bought by the same customers, **regardless of** customer clusters
- ❖ Co-clustering (2 variables):
 - Group customers and products **simultaneously**
 - Two customers are similar if they buy similar products
 - Two products are similar if they are bought by similar customers

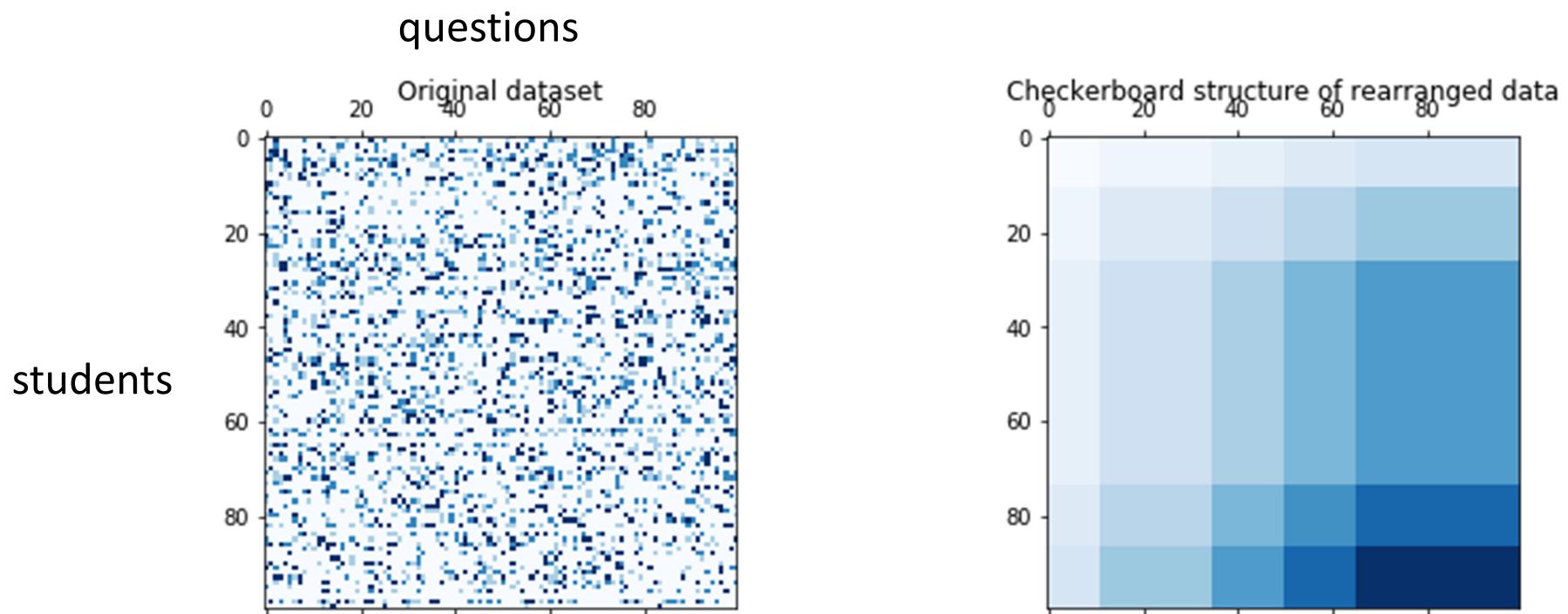
Co-clustering: Example

- ❖ Clustering two variables **simultaneously**



Co-clustering: Another Example

- ❖ Detect whether students copy exam from one another

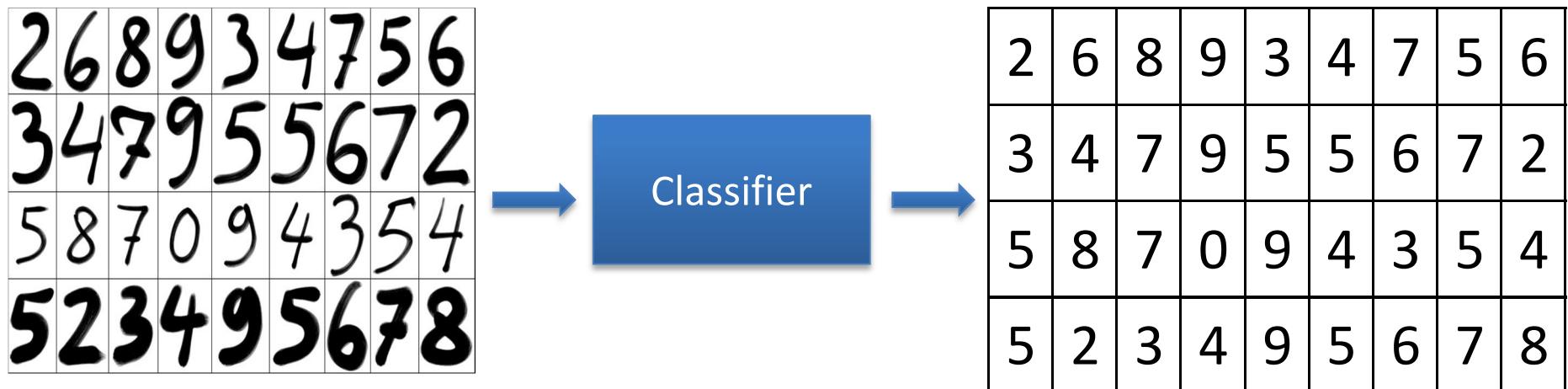


II. Classification

- ❖ **Task:** predict the **label/class** of data
- ❖ **How it works** (supervised learning):
 - Learn from a collection of pre-labeled data (training set)
 - Each record contains a set of attributes or features
 - Each record is assigned with a label/class
 - Build a model for calculating the class as a function of the attributes.
 - Evaluate the model on previously unseen records (test set)
- ❖ **Goal:** increase the accuracy of test set
- ❖ **Applications:**
 - Pattern Recognition
 - Time Series Prediction
 - Signal Processing
 - Anomaly Detection

Classification: Other Applications

- ❖ Handwritten digit recognition:



- ❖ Email spam recognition:



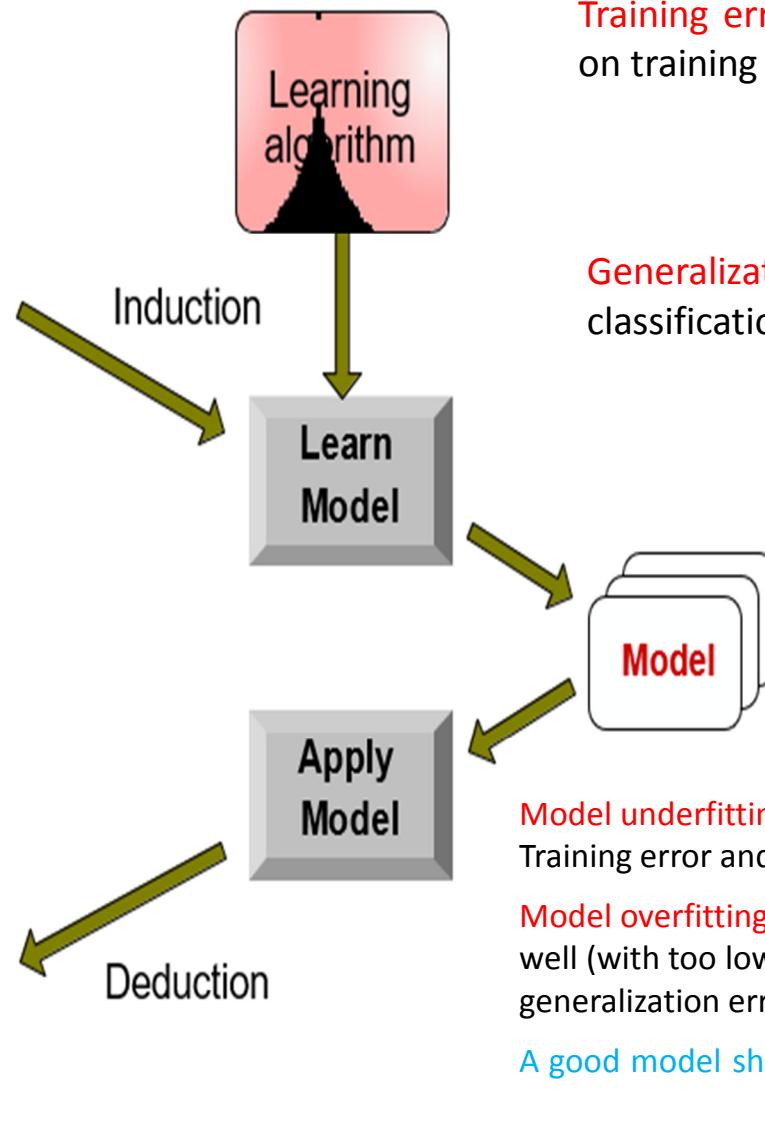
Classification Pipeline

Tid	Attrib1	Attrib2	Attrib3	Class
1	Yes	Large	125K	No
2	No	Medium	100K	No
3	No	Small	70K	No
4	Yes	Medium	120K	No
5	No	Large	95K	Yes
6	No	Medium	60K	No
7	Yes	Large	220K	No
8	No	Small	85K	Yes
9	No	Medium	75K	No
10	No	Small	90K	Yes

Training Set

Tid	Attrib1	Attrib2	Attrib3	Class
11	No	Small	55K	?
12	Yes	Medium	80K	?
13	Yes	Large	110K	?
14	No	Small	95K	?
15	No	Large	67K	?

Test Set



Training errors: number of classification errors on training records.

Generalization (test) errors: number of classification errors on test records.

Model underfitting: model is too simple such that both Training error and Generalization error are high.

Model overfitting: A model that fits the training data too well (with too low training errors) may have a higher generalization error than a model with higher training error.

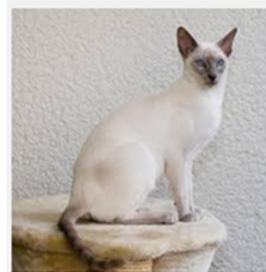
A good model should have low errors of both types.

Classification: Methods

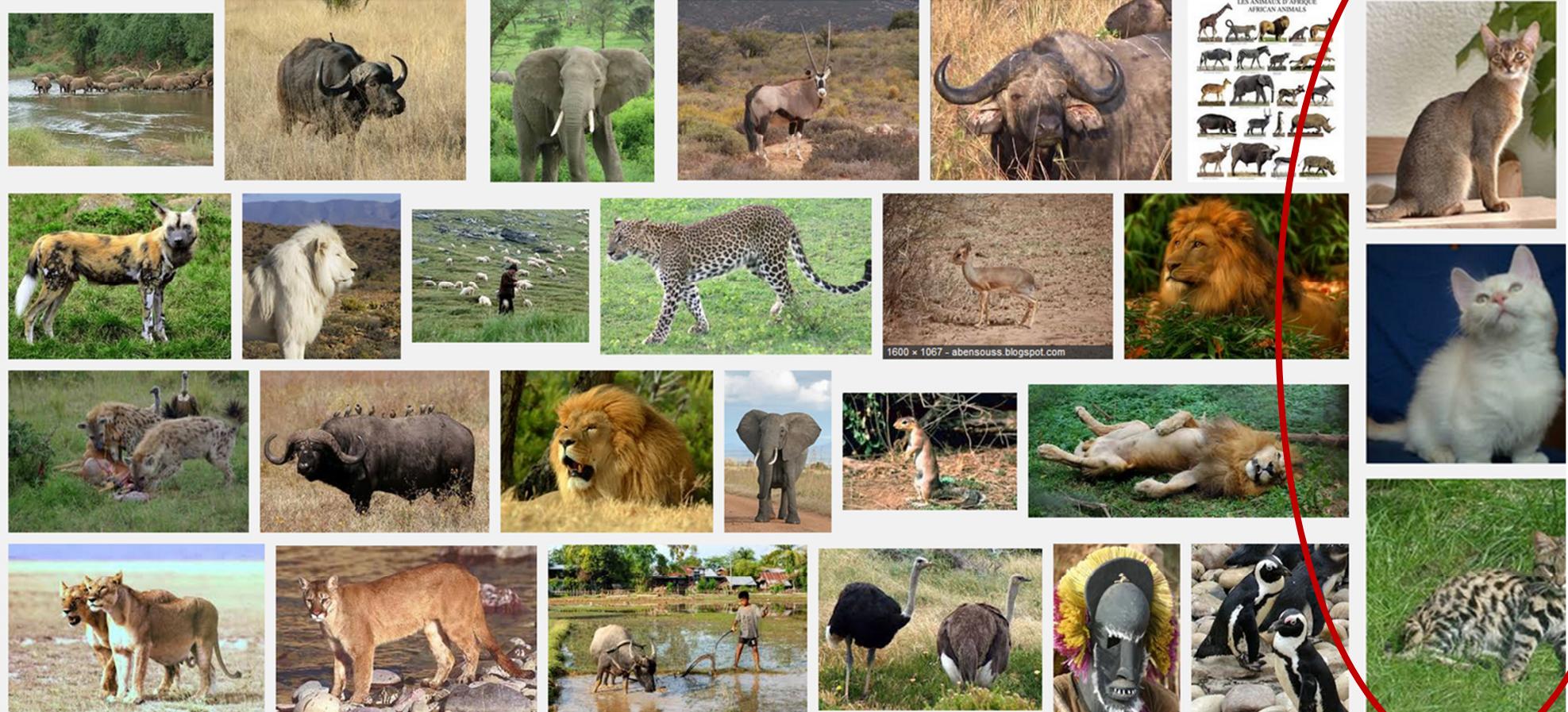
1. K Nearest Neighbors (k-NN)
2. Decision Tree
3. Naïve Bayes
4. Linear Classifier: Support vector machine
5. Logistic Classifier: SVM with kernels
6. Supervised neural networks
7. Ensemble method

1. k-NN: Example

Find the label for a query item
from k closest matches
in a labeled dataset



With k=3, majority vote
should be "cat"



k-NN

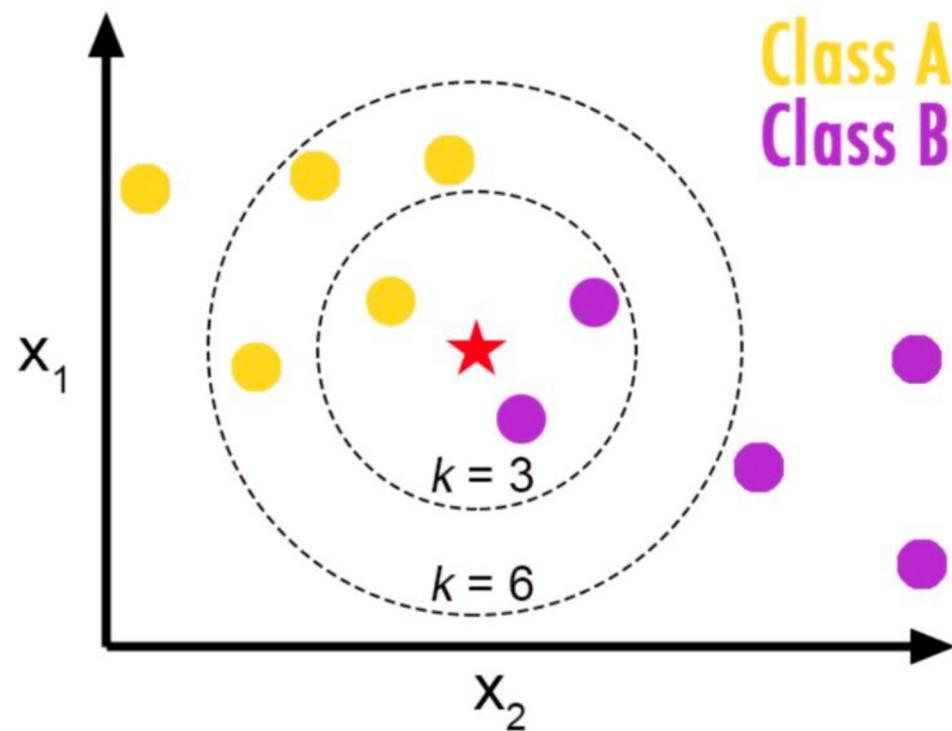
- ❖ **Training Time:** memorize all data and labels
- ❖ **Testing Time:** predict the label of a new data by take majority vote from K closest points
 - Can also use a weighted vote of the neighbors.
 - Can use different similarity metrics:
 - Euclidean distance: simplest, fast to compute, geometric space
$$d(x, y) = \|x - y\|$$
 - Cosine distance: good for documents, images
$$d(x, y) = 1 - \frac{x \cdot y}{\|x\| \cdot \|y\|}$$
 - Jaccard distance: for set data
$$d(X, Y) = 1 - \frac{|X \cap Y|}{|X \cup Y|}$$
 - Hamming distance: for string data
$$d(x, y) = \sum_{i=1}^n 1_{x_i \neq y_i}$$

Distance metrics

- ❖ Given two vectors: $x = (1,0,1)$, $y = (0,1,1)$
 - $X = \{1,3\}$, $Y = \{2,3\}$ (set of index numbers of x,y having value=1)
 - Euclidean distance: $d(x, y) = \|(1, -1, 0)\| = \sqrt{2}$
 - Cosine distance: $d(x, y) = 1 - \frac{1}{\sqrt{2}\sqrt{2}} = 1/2$
 - Jaccard distance: $d(X, Y) = 1 - \frac{|\{1,3\} \cap \{2,3\}|}{|\{1,3\} \cup \{2,3\}|} = 2/3$
 - Hamming distance: $d(x, y) = 1 + 1 = 2$

k-NN

- ❖ Choosing a K will affect what class a new point is assigned to:

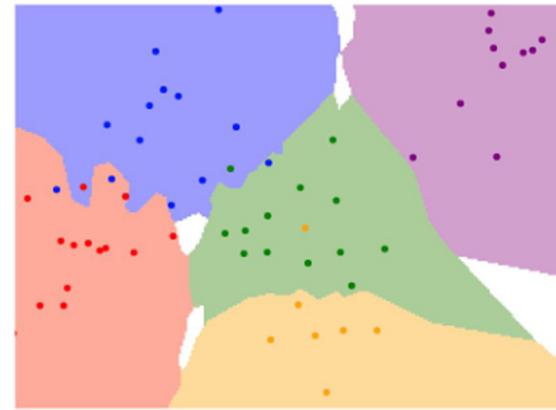


k-NN: Example

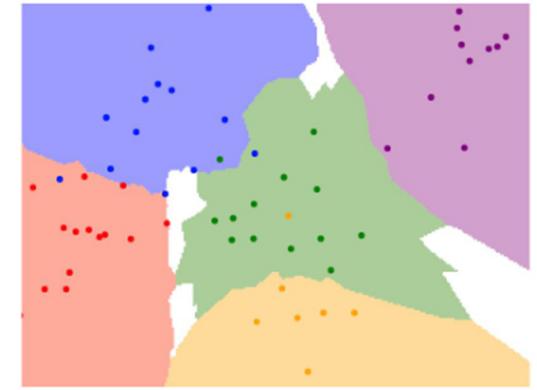
- ❖ We can draw a **classification map** for each k:
 - Any new test point fallen into a colored area takes the color (label) of that area.



k=1



k=3



k=5

K-NN: pros and cons

❖ Pros:

- Very simple
- Training is trivial
- Works with any number of classes
- Easy to add more data
- Few parameters: k, distance metric

❖ Cons:

- High **computation cost** for testing (worse for large data sets)
- Not good with high dimensional data
 - The curse of dimensionality: similar data in low-dimensional space might be far away from each other in high-dimensional space
- Categorical features don't work well

2. Decision Tree

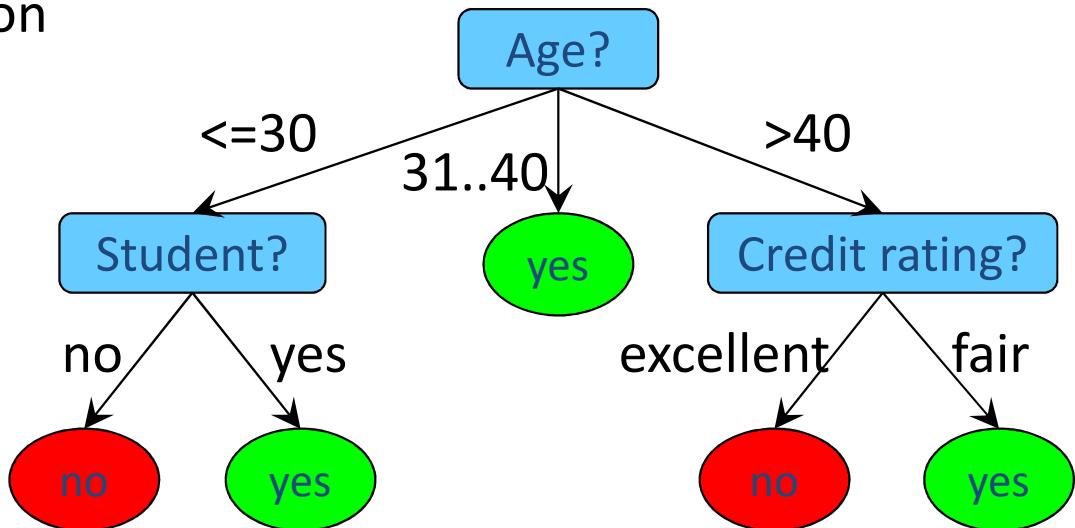
- ❖ **Example:**
- ❖ **Model:** flow-chart-like tree structures

- Nodes are tests on a single attribute
- Branches are attribute values
- Leaves are marked with class labels

- ❖ **Training time:** optimize classification accuracy (decision tree induction)

- NP-hard
- Heuristic: top-down tree construction + pruning

age	income	student	credit_rating	buys_computer
<=30	high	no	fair	no
<=30	high	no	excellent	no
31..40	high	no	fair	yes
>40	medium	no	fair	yes
>40	low	yes	fair	yes
>40	low	yes	excellent	no
31..40	low	yes	excellent	yes
<=30	medium	no	fair	no
<=30	low	yes	fair	yes
>40	medium	yes	fair	yes
<=30	medium	yes	excellent	yes
31..40	medium	no	excellent	yes
31..40	high	yes	fair	yes
>40	medium	no	excellent	no



Decision Tree: Pros and Cons

❖ Pros:

- Easy to interpret **visually**
- Work well with qualitative (categorical) features

❖ Cons:

- Prone to over-fitting
- **Non-incremental**: need to be re-trained from scratch if new training data becomes available

3. Naïve Bayes Classifier

- ❖ Given a new instance, estimate for each possible value of the class attribute its probability.
 - Choose the value with the maximum probability.

Tid	Refund	Marital Status	Taxable Income	Cheat
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

Refund	Marital Status	Taxable Income	Cheat
No	Married	80K	?

1. Estimate conditional probabilities
 $P1 = P(\text{cheat}=\text{yes} | \text{Ref}=\text{No}, \text{Ms}=\text{Married}, \text{Income}=80)$
 $P2 = P(\text{cheat}=\text{No} | \text{Ref}=\text{No}, \text{Ms}=\text{Married}, \text{Income}=80)$

If $P1 > P2$, classify as cheat, otherwise no cheat

How to estimate the probabilities?

Probability: Revisit

- ❖ If events A and B are independent, then $P(B|A) = P(B)$, and the probability that event A and B both occur, denoted $P(A,B)$, is equal to $P(A)P(B)$.
- ❖ Bayes' Theorem:

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

$$P(B|A)P(A) = P(A|B)P(B) = P(A, B)$$

- ❖ Chain Rule: permits the calculation of any member of the joint distribution of a set of random variables using only conditional probabilities.

➤ e.g. $P(a_4, a_3, a_2, a_1) = P(a_4|a_3, a_2, a_1) \cdot P(a_3|a_2, a_1) \cdot P(a_2|a_1) \cdot P(a_1)$

Naïve Bayes Classifier

❖ Model:

- Consider each attribute and class label as **random variables**
- Give a record with attributes $A = (a_1, \dots, a_n)$
 - Goal is to predict class C
 - Specifically, we want to find the value of C that maximizes:
 $P(C|a_1, \dots, a_n)$
 - e.g. estimate $P(\text{cheat=yes} | \text{Refund=No}, \text{MStatus=Married}, \text{Income=80K})$ or shorthand $P(\text{yes} | (\text{No, Married, 80K}))$

Naïve Bayes Classifier

❖ **Solution:** use Bayes theorem $P(C|A) = \frac{P(A|C)P(C)}{P(A)}$

- No need to estimate $P(A)$ for comparisons of $P(C = c_1|A)$ and $P(C = c_2|A)$
- $P(C = c_k)$ can be estimated as the fraction of training records taking class value c_k

$$P(c_k) = \frac{N_k}{N}$$

- Estimate $P(A|c_k)$ is hard: assume the features a_1, \dots, a_n are independent given class C .

$$P(A|c_k) = \prod_{i=1}^n P(a_i|c_k)$$

↑
 $\frac{N_{ik}}{N_k}$

N : Total number of data

N_k : Total number of records in class k

N_{ik} : Total number of records with value a_i for feature i in class k

Zero conditional probability

- ❖ If no sample contains the feature value
 - We face a **zero conditional probability** problem during testing:
$$\prod_{i=1}^N P(a_i|c_k) = 0 \text{ if any } P(a_i|c_k) = 0$$
 - Thus, we have to **re-estimate** class conditional probabilities as:
$$P(a_i|c_k) = \frac{N_{ik} + mp}{N_k + m}$$
 - $m \geq 1$: weight to prior (number of “virtual” examples, usually up to 1% number of training samples with class c_k)
 - p : prior estimate (usually, $p = 1/t$ for t possible values of a_i)

Naïve Bayes Classifier: Example

❖ E.g. compute

$$P(C|No, Married, 80K) \sim P(No, Married, 80K|C) P(C)$$

- $P(C=Yes) = 3/10, P(C=No) = 7/10$
- $P(No|C=Yes) = 3/3, P(No|C=No) = 4/7$
- $P(Married|C=Yes) = 0, P(Married|C=No) = 3/7$
→ Have to use zero conditional probability
 - $P(Married | C=Yes) = (0+1*1/3)(3+1) = 1/12$
- $P(80K | C=Yes) = 0, P(80K | C=No) = 0$
→ Have to use zero conditional probability
 - $P(80K | C=Yes) = (0+1*1/10)/(3+1) = 1/40$
 - $P(80K | C=No) = (0+1*1/10)/(7+1) = 1/80$

❖ Put it all:

- $P(C=Yes|No, Married, 80) = 3/10 * 3/3 * 1/12 * 1/40 = 0.000625$
- $P(C=No|No, Married, 80) = 7/10 * 4/7 * 3/7 * 1/80 = 0.00214$
- So the new sample is not likely to cheat

Tid	Refund	Marital Status	Taxable Income	Cheat
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

Refund	Marital Status	Taxable Income	Cheat
No	Married	80K	?

Naïve Bayes Classifier: Pros and Cons

❖ Pros:

- Computationally fast
- Simple to implement
- Works well with high dimensions

❖ Cons:

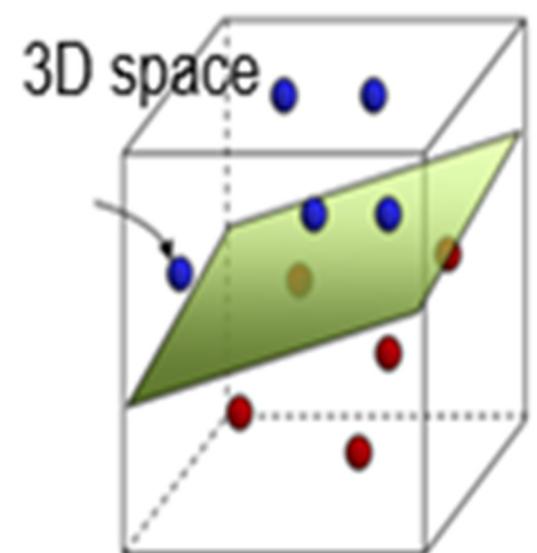
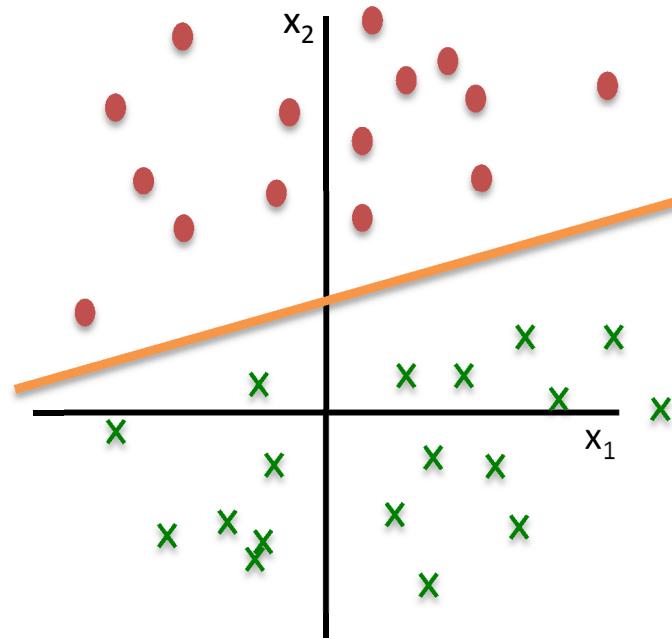
- Relies on independence assumption
- Not work well with categorical inputs (zero conditional probability issue)

4. Linear Classification: Support Vector Machine

- ❖ We want to draw a line to separate two classes

Linearly separable case!

An Ideal case



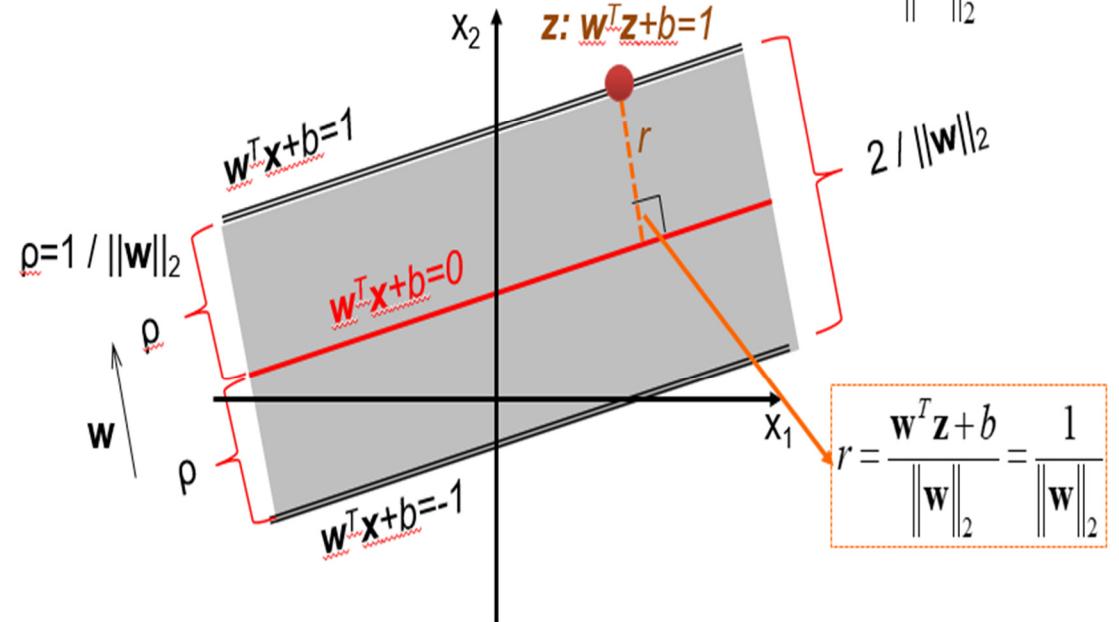
SVM: Core Idea

- ❖ **Model:** we focus on two parallel hyperplanes:

$$\begin{cases} \mathbf{w}^T \mathbf{x} + b = 1 \\ \mathbf{w}^T \mathbf{x} + b = -1 \end{cases}$$

Geometrically, the distance between these two planes is

$$\frac{2}{\|\mathbf{w}\|_2}$$



SVM - Separation Margin

- Given two parallel hyperplanes below, we separate two classes of data by preventing samples from falling into the margin:

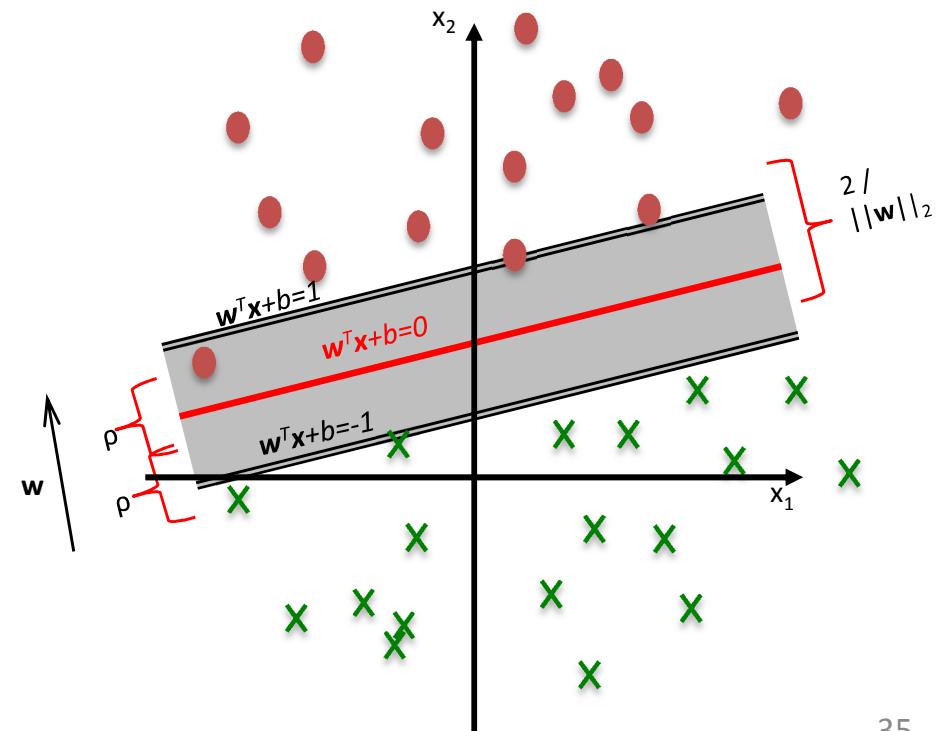
$$\begin{cases} \mathbf{w}^T \mathbf{x} + b \geq 1, & \text{if } y=1, \\ \mathbf{w}^T \mathbf{x} + b \leq -1, & \text{if } y=-1. \end{cases}$$

equivalent expression
 ⇌

$$y(\mathbf{w}^T \mathbf{x} + b) \geq 1$$

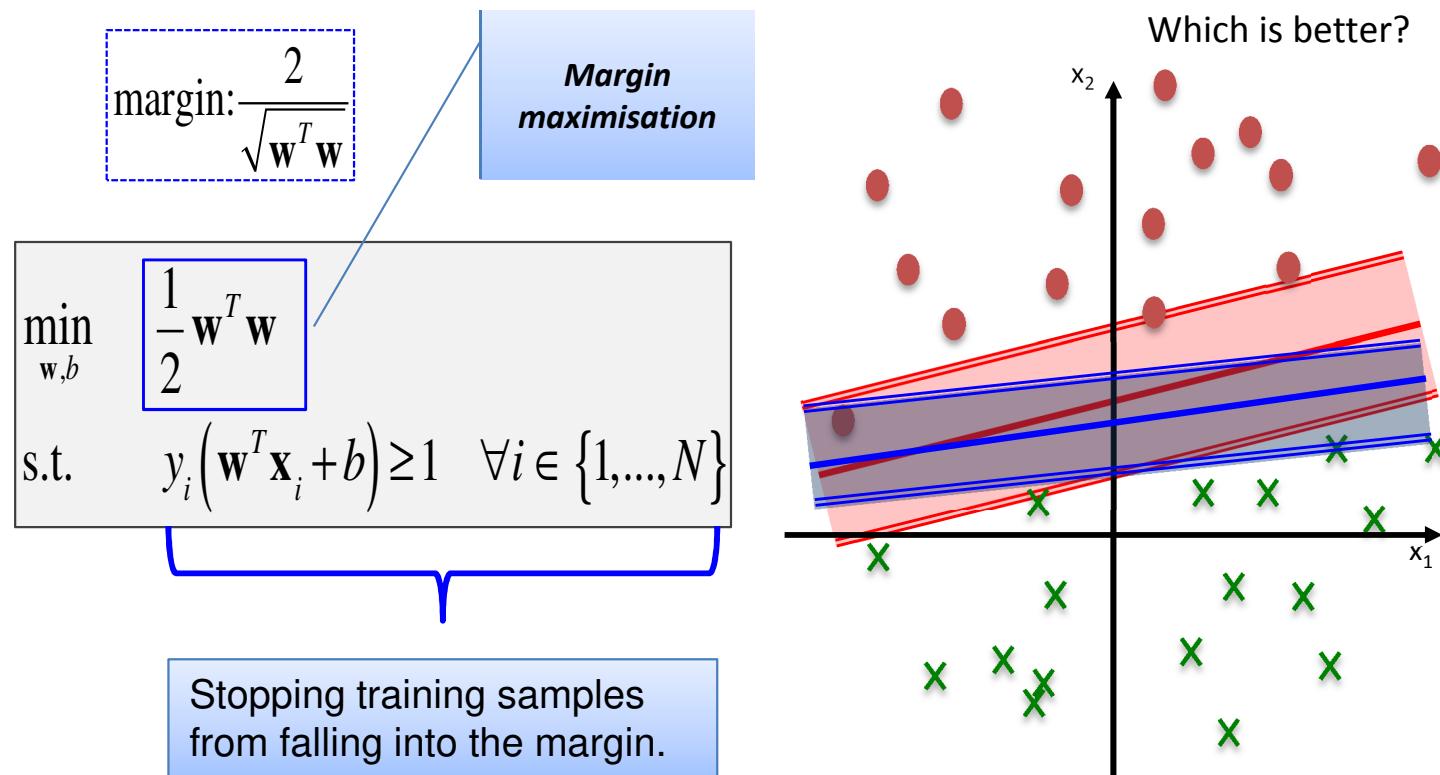
- The region bounded by these two hyperplanes is called the separation “margin”, given by

$$\rho = \frac{2}{\|\mathbf{w}\|_2} = \frac{2}{\sqrt{\mathbf{w}^T \mathbf{w}}}$$



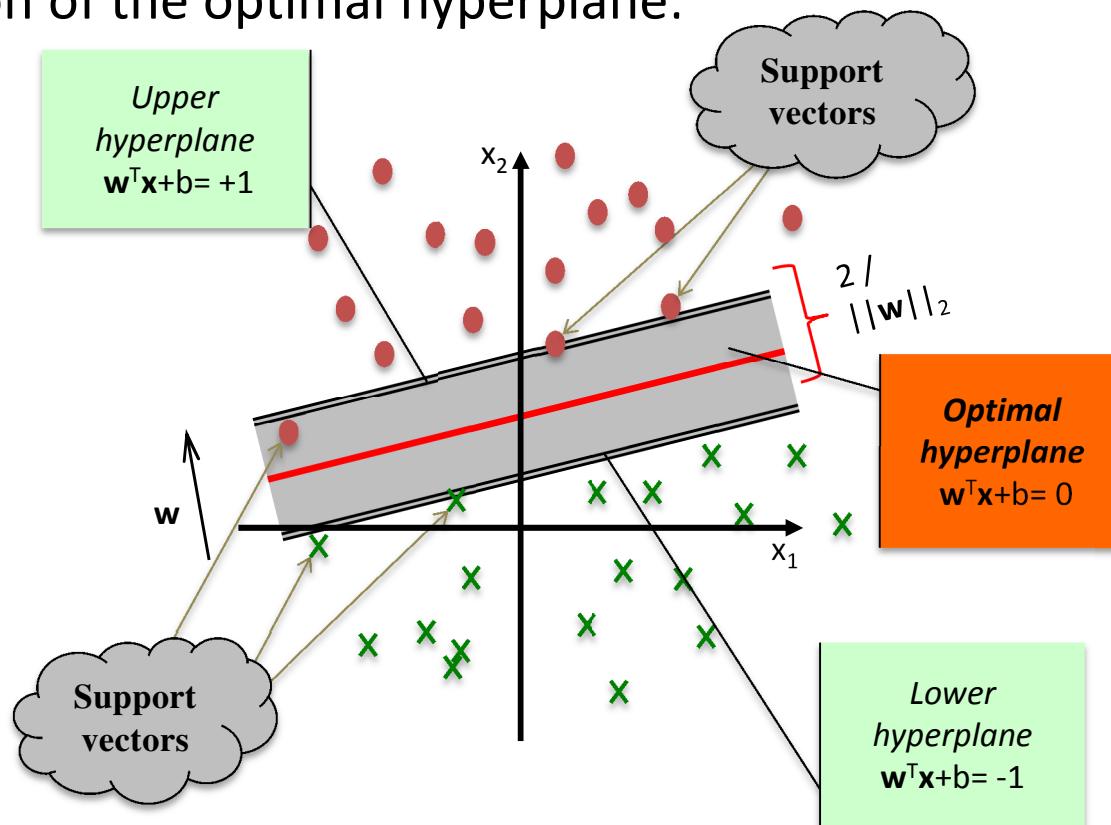
SVM: Objective function

- ❖ **Goal:** The aim of SVM is simply to find an optimal hyperplane to separate two-classes of data points with the widest margin.
- ❖ This results in the following constrained optimisation:



SVM - Support Vectors

- ❖ Support vectors: data points that satisfy $y(\mathbf{w}^T \mathbf{x} + b) = 1$
- ❖ These points are the most difficult to classify and are very important for the location of the optimal hyperplane:



SVM - Training (Optional)

- ❖ **SVM training:** the process of solving the following constrained optimisation problem (aka hard-margin SVM):

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} \\ \text{s.t.} \quad & y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1 \quad \forall i \in \{1, \dots, N\} \end{aligned}$$

Lagrange

Dual problem

$$\begin{cases} \max_{\lambda \in \mathbb{R}^N} \sum_{i=1}^N \lambda_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j \\ \text{s.t.} \quad \sum_{i=1}^N \lambda_i y_i = 0 \\ \lambda_i \geq 0 \end{cases}$$

- ❖ **Solution:** quadratic programming https://en.wikipedia.org/wiki/Quadratic_programming

- Following calculus and optimisation theory, we use Lagrange multipliers $\lambda_i \geq 0$ and formulate the Lagrangean function:

$$L(b, \mathbf{w}, \lambda) = \frac{1}{2} \mathbf{w}^T \mathbf{w} + \sum_{i=1}^N \lambda_i [1 - y_i (\mathbf{w}^T \mathbf{x}_i + b)]$$

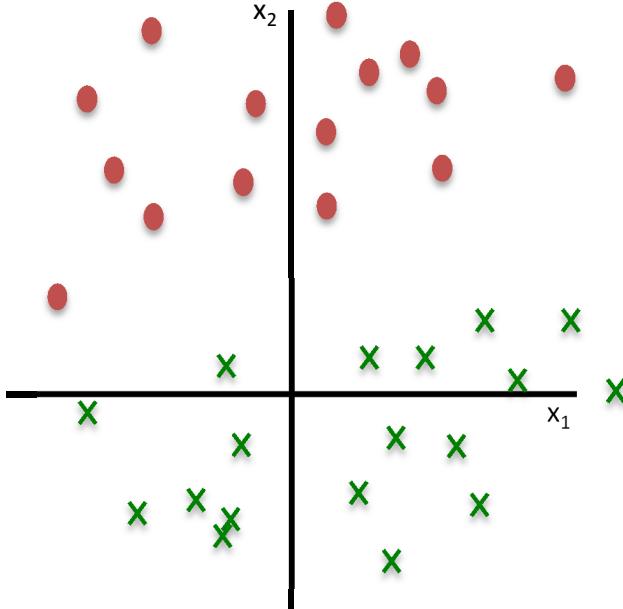
$$\begin{aligned} \frac{\partial L}{\partial b} = 0 &\Rightarrow \sum_{i=1}^N \lambda_i y_i = 0 \\ \frac{\partial L}{\partial \mathbf{w}} = 0 &\Rightarrow \sum_{i=1}^N \lambda_i y_i \mathbf{x}_i = \mathbf{w} \end{aligned}$$

- By differentiating $L()$ and setting the gradient to zero, we can express the Lagrangean function using only the multipliers (called a dual problem).

$$L(\lambda) = \sum_{i=1}^N \lambda_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$$

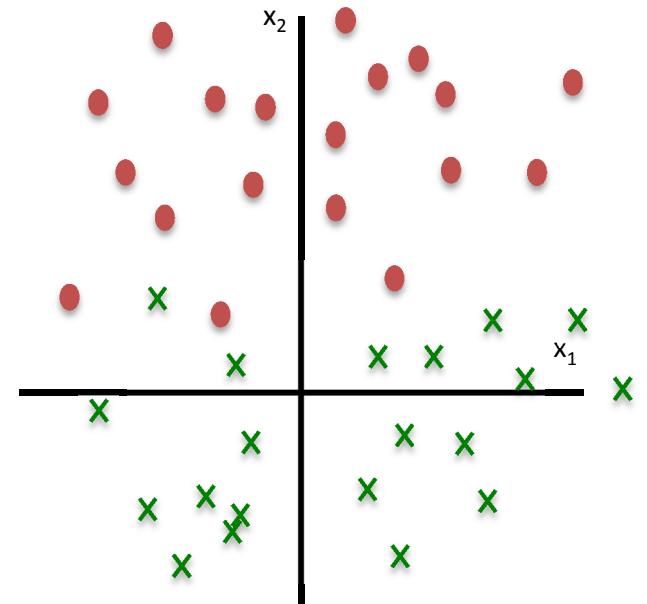
SVM: Soft Margin

So far, we work on simple cases like this:



separable data patterns

What if the data points look like this?



non-separable data patterns

In practice, no datasets are ideally linearly separable. This means that some data points are bound to be misclassified with a linear hyperplane.

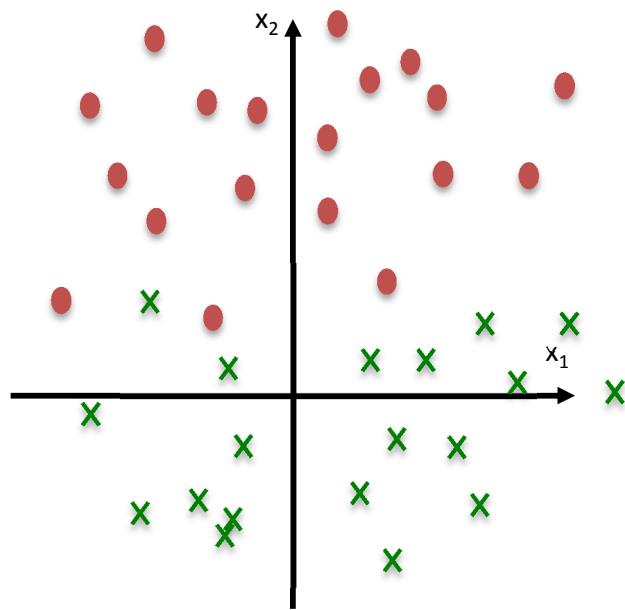
SVM with softmax margins

$$\begin{aligned} \min_{(\mathbf{w}, b) \in \Re^{d+1}} & \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^N \xi_i \\ \text{s.t.} & \left. \begin{aligned} y_i (\mathbf{w}^T \mathbf{x}_i + b) &\geq 1 - \xi_i \\ \xi_i &\geq 0 \end{aligned} \right\} \forall i \in \{1, \dots, N\} \end{aligned}$$

$C \geq 0$ is a user defined parameter, which controls the **regularisation**. This is the trade-off between complexity and non-separable patterns.

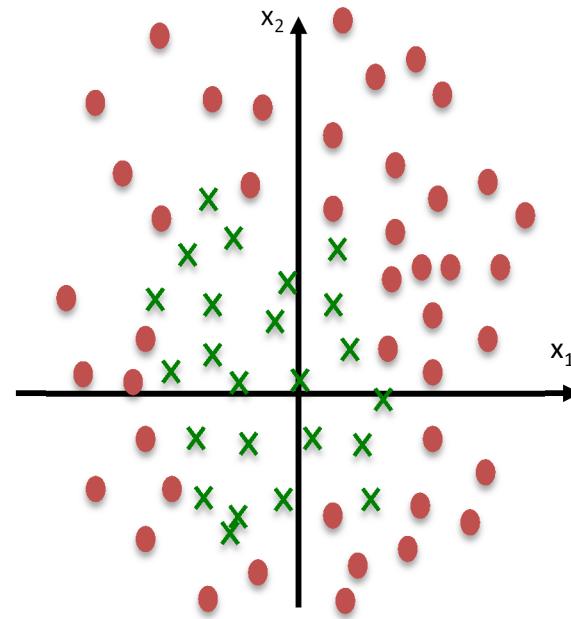
5. Non-linear Classifier

So far, we can handle linear cases like this:



linear data patterns

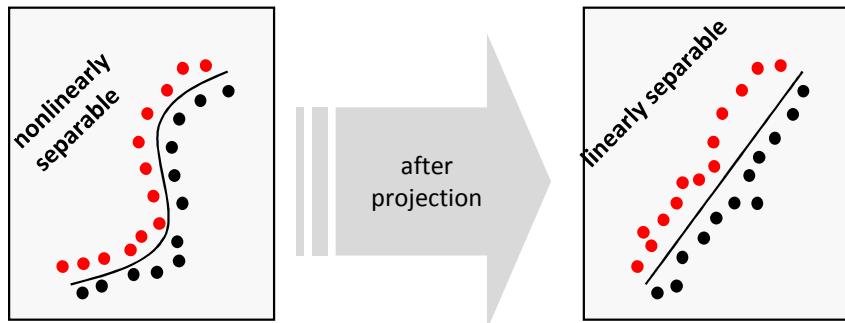
What if the data points look like this?



non-linear data patterns

SVM - Kernel Methods

- ❖ A method to handle nonlinear data patterns.
- ❖ Similar to linear basis function model: we will **project each data point to a new feature space** in order to make the patterns linearly separable in that space.



Kernel SVM demos:

<https://www.youtube.com/watch?v=3liCbRZPrZA>

<https://www.youtube.com/watch?v=ndNE8he7Nnk>

- ❖ Commonly used Kernel functions

Kernel Name	Expression $K(\mathbf{x}, \mathbf{y})$	Comments
Linear	$\mathbf{x}^T \mathbf{y}$	No parameter
Polynomial	$(\mathbf{x}^T \mathbf{y} + 1)^p$	p is a user defined parameter
Gaussian, also called radial basis function (RBF)	$\exp\left(-\ \mathbf{x} - \mathbf{y}\ _2^2 / (2\sigma^2)\right)$	σ is the user defined width parameter

SVM: Pros and Cons

❖ Pros:

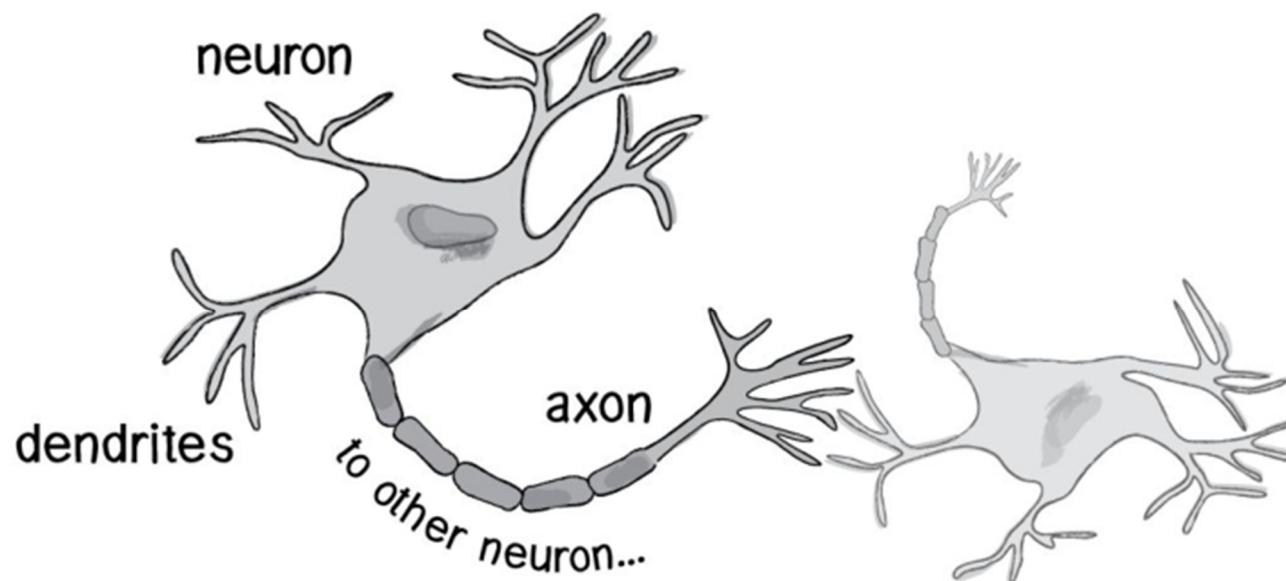
- Handle high dimensional data well
- Work with both linear and non-linear boundary depending on the kernel used

❖ Cons:

- Susceptible to over-fitting depending on kernel → Different kernels have different uses (No clear winner)

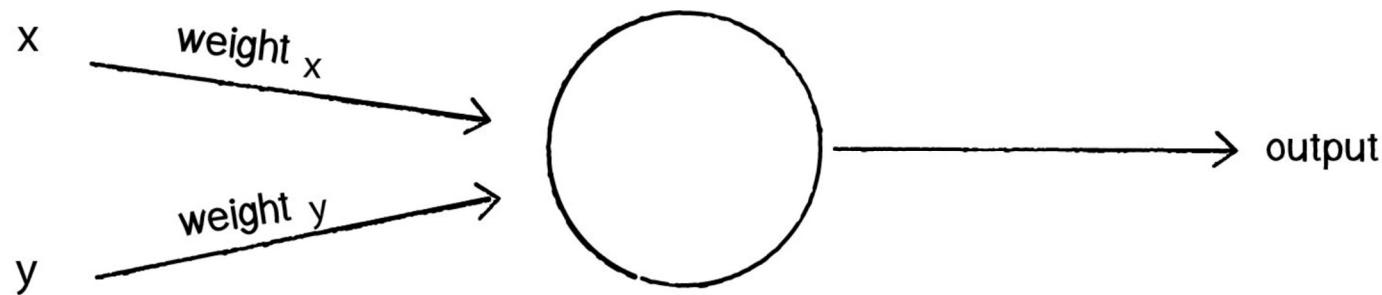
6. Neural Networks

- ❖ Neural networks are **modeled after** biological neural networks and attempt to allow computers to learn in a similar manner to humans
- ❖ The human brain has interconnected neurons with dendrites that receive inputs, and then based on those inputs, produce an electrical signal output through the axon



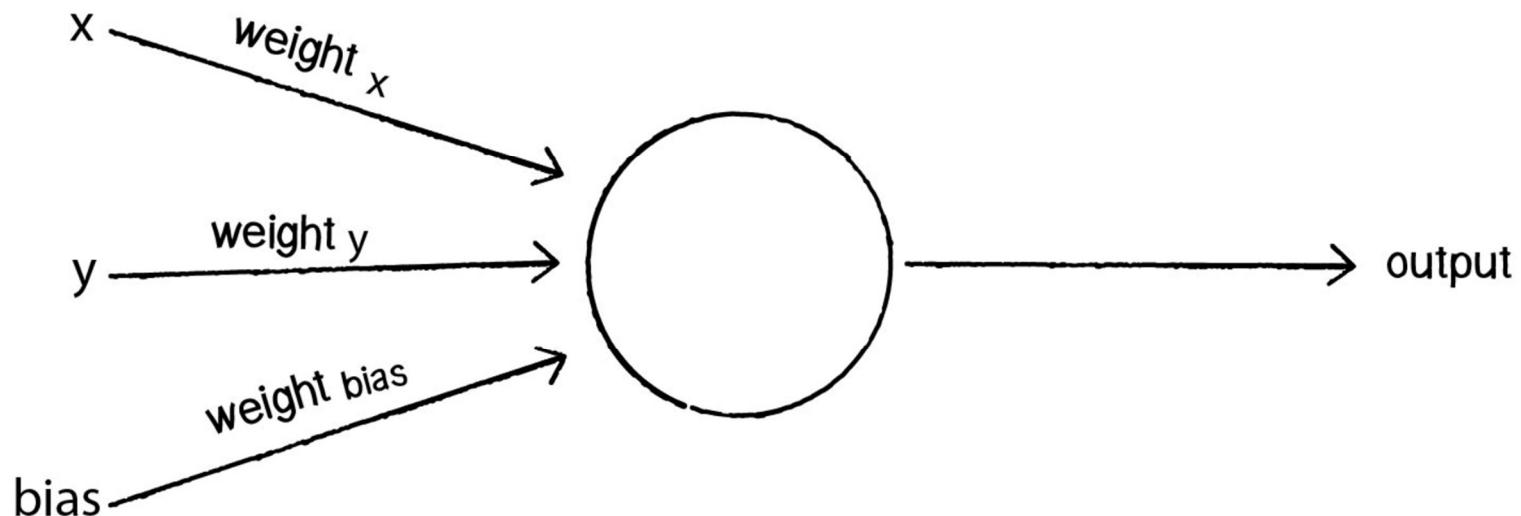
6. Neural Networks

- ❖ **Perceptron** (the simplest neural network): follows the “**feed-forward**” model, meaning inputs are sent to the neuron, are processed, are result in an output.
 1. Receive inputs
 2. Weight inputs
 3. Sum Inputs
 4. Generate Output
- ❖ E.g. $f(x,y) = \text{weight}_x * x + \text{weight}_y * y$
- ❖ **Learning** = find the “best” weights for a set of inputs



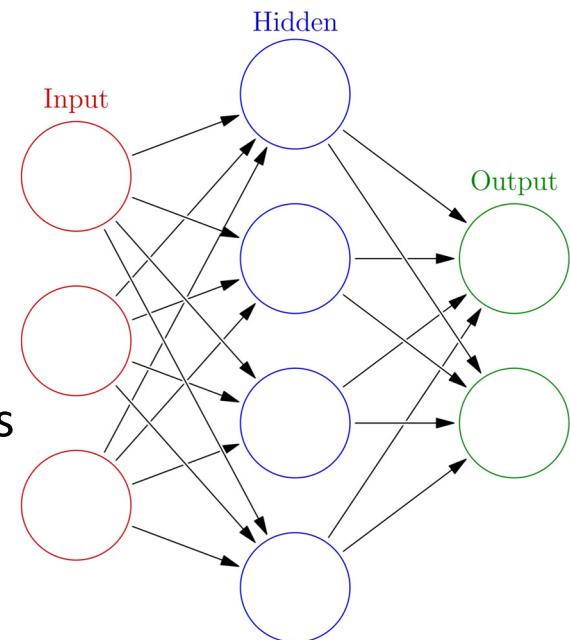
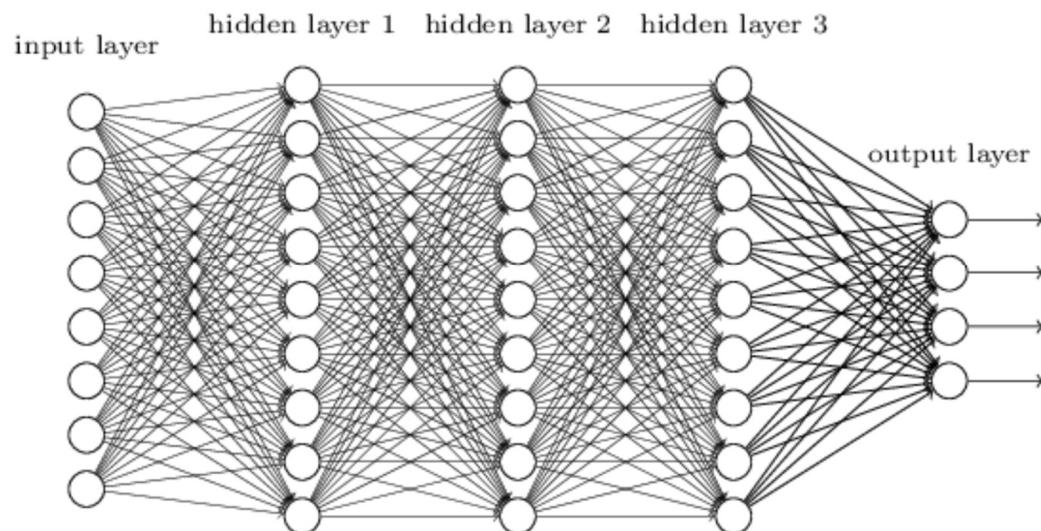
Neural Networks

- ❖ Sometimes we use a **bias factor** to avoid zero issue:
 - E.g. if both inputs are zero, then weights cannot be learned



Neural networks

- ❖ In **general** neural networks, there are 3 types of layers:
 - Input layer
 - Hidden layer
 - Output layer
- ❖ Where does the term "**Deep Learning**" come from?
 - A neural network with many hidden layers
 - Microsoft won vision recognition contest with 152 layers



Neural networks

- ❖ **General learning process:**

1. Provide the inputs for which there is a known answer
2. Ask the network to guess an answer
3. Compute the error (how far off the correct answer?)
4. Adjust all the weights according to the error
5. Return to step 1 and repeat until we are satisfied with the error

- ❖ **More details:** see Deep Learning courses

Neural networks: pros and cons

❖ Pros:

- State-of-the-art model with very high accuracy
- Learning process can be parallelized and incremental

❖ Cons:

- Computationally expensive (need special hardware such as GPU)
- Hard to choose parameters (e.g. number of layers, types of layers, width of layers)

7. Classification: Ensemble Method

Are like **crowdsourced machine learning algorithms**:

- ❖ Take a collection of simple or *weak* learners
- ❖ Combine their results to make a single, better learner

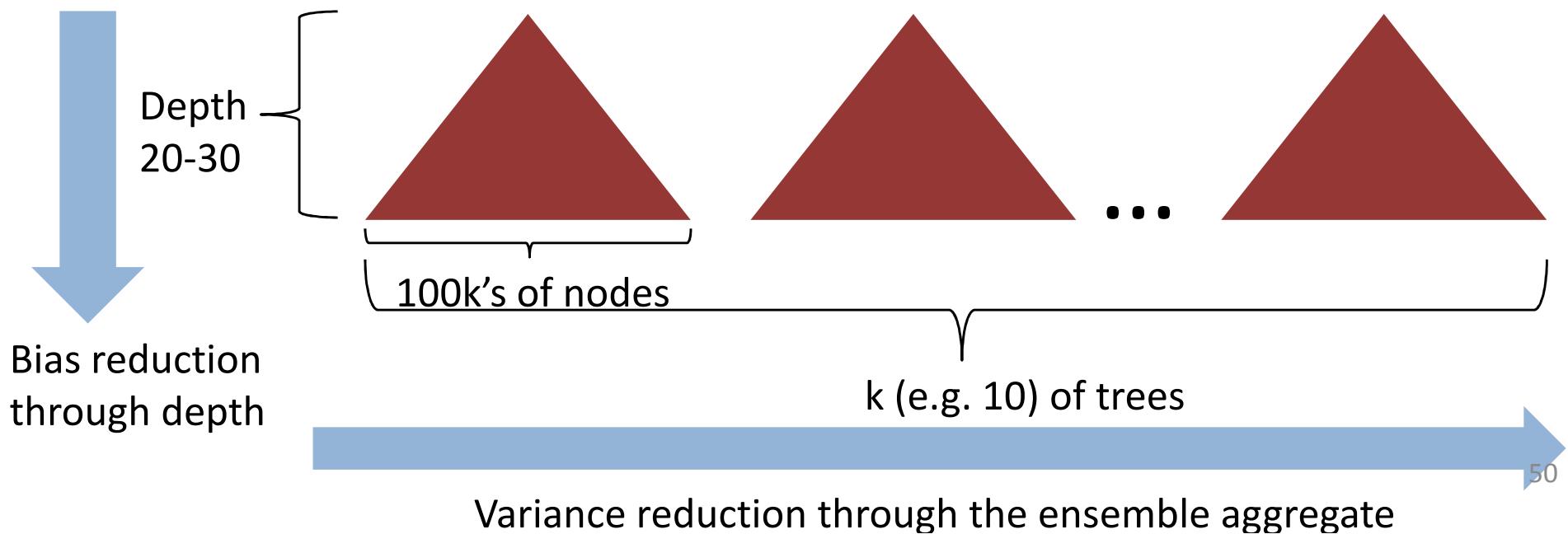
Types:

- ❖ **Bagging**: train learners in parallel on different samples of the data, then combine by voting (discrete output) or by averaging (continuous output).
- ❖ **Stacking**: combine model outputs using a second-stage learner like linear regression.
- ❖ **Boosting**: train learners on the filtered/weighted/... output of other learners.

Ensemble Method: Random Forest (Optional)

Combine multiple decision trees:

- Draw K bootstrap **samples** of size N from original dataset
- Grow a decision tree for each sample, by selecting a **random set of m out of p features** at each node, and choosing the best feature to split on. (typically $m = \sqrt{p}$ but can be smaller)
- Aggregate the predictions of the trees (most popular vote) to produce the final class.



Ensemble Method: Pros and Cons

❖ Pros:

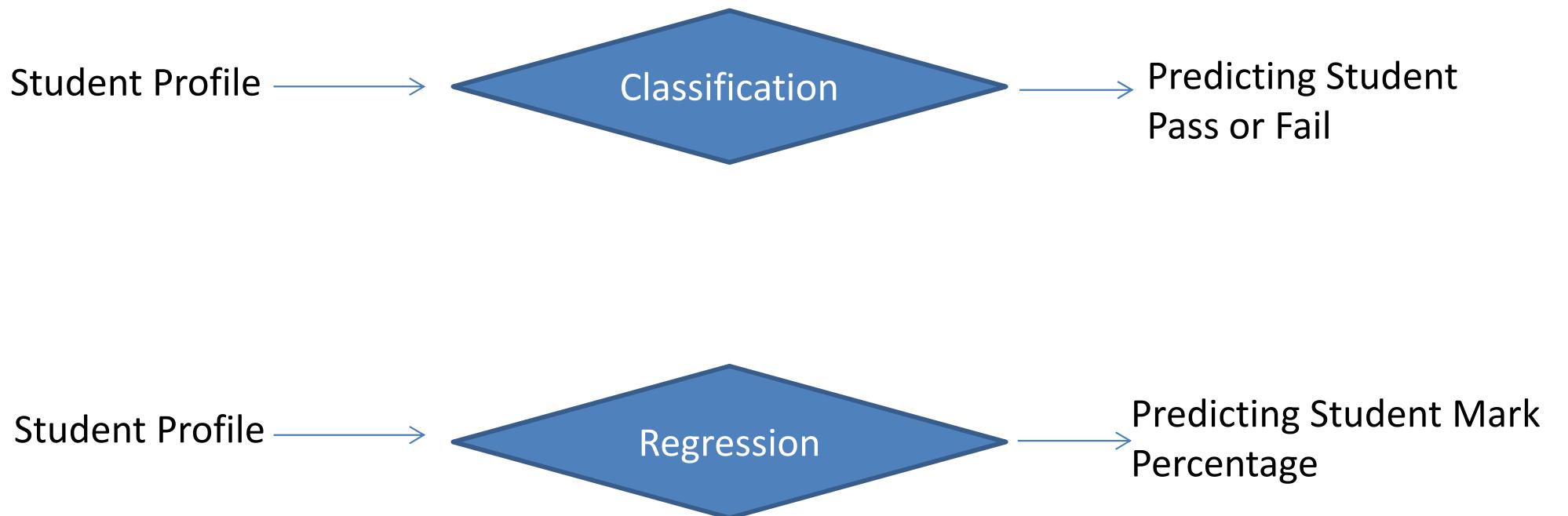
- Effective for dense data (a few thousand features)
- Parallelizes easily

❖ Cons:

- Not easy to implement (how many classifiers are enough?)
- Needs many passes over the data – at least the max depth of the trees.

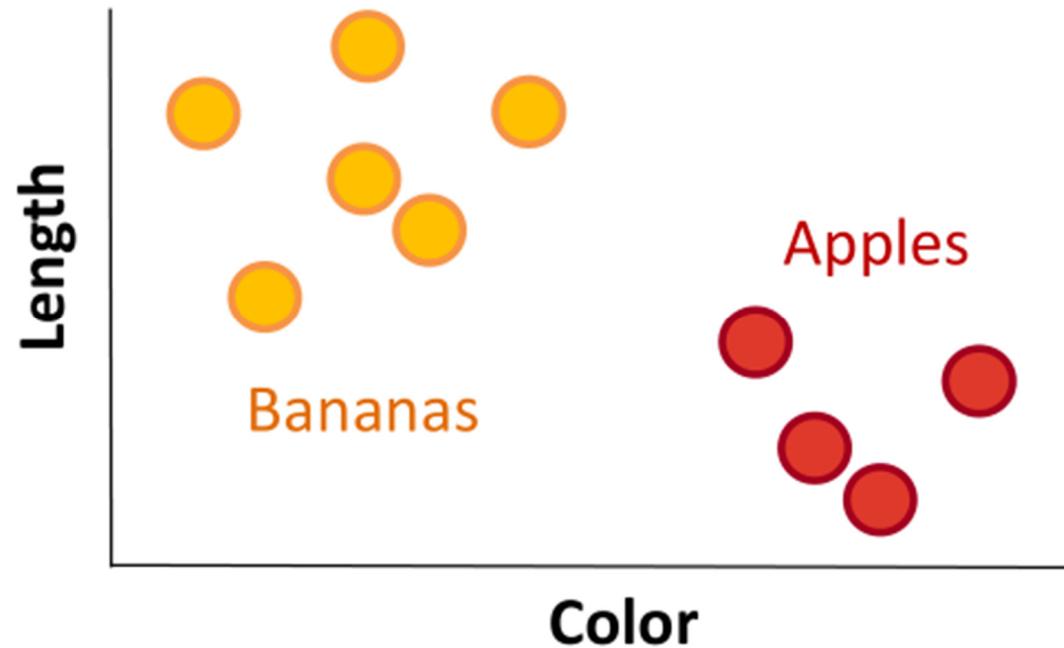
III. Regression

❖ Classification vs. Regression:



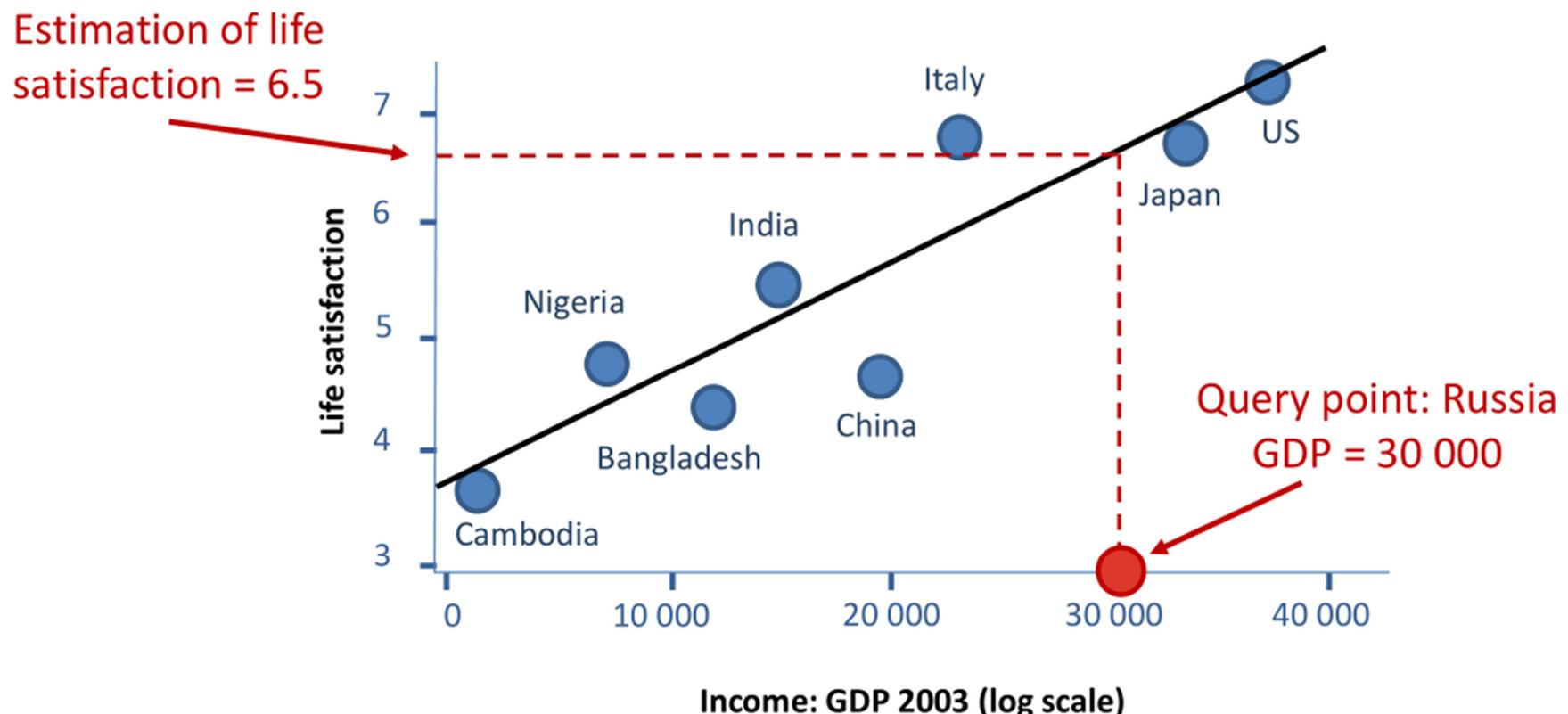
Classification: Recap

- ❖ Maps N -dimensional input $x \in R^N$ to **discrete** values y
- ❖ e.g. $x = [\text{Length}, \text{Color}]$, $y = \{\text{"Banana"}, \text{"Apple"}\}$



Regression: Definition

- ❖ Maps N -dimensional input $x \in R^N$ to **continuous** values $y \in R$
- ❖ e.g. $x = [\text{Income (GDP)}]$, $y = \text{Continuous value of life satisfaction}$



Regression: Applications

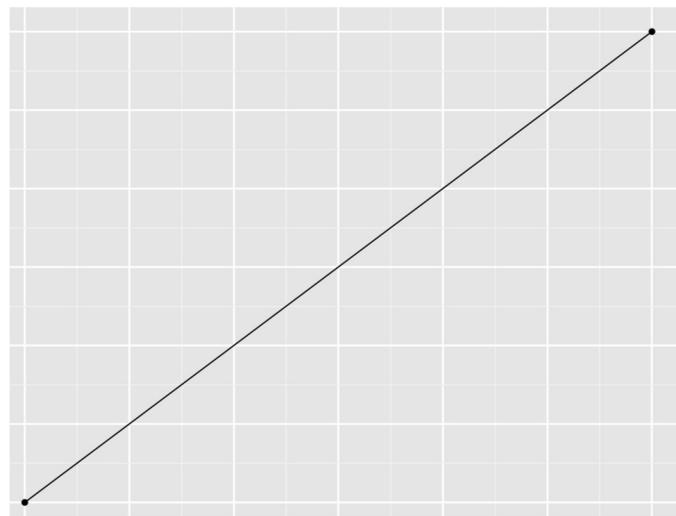
- ❖ **Predictive analytics:** forecasting future opportunities and risks
 - Forecast the number of shoppers who will pass in front of a particular billboard
 - Use that model to estimate the price for an advertisement
- ❖ **Operation Efficiency:** optimize business processes
 - Call center: correlation between wait times of callers vs. number of complaints
 - Use that model to reduce the number of complaints to some degree
- ❖ **Support Decisions:** understand customer behavior to make informed business decisions
- ❖ **Correcting Errors:** verify assumptions and adjust business strategies
- ❖ **New Insights:** Uncover patterns that were previously unnoticed
 - E.g. Market demand increases on certain days of the week or certain times of the year

Regression Models

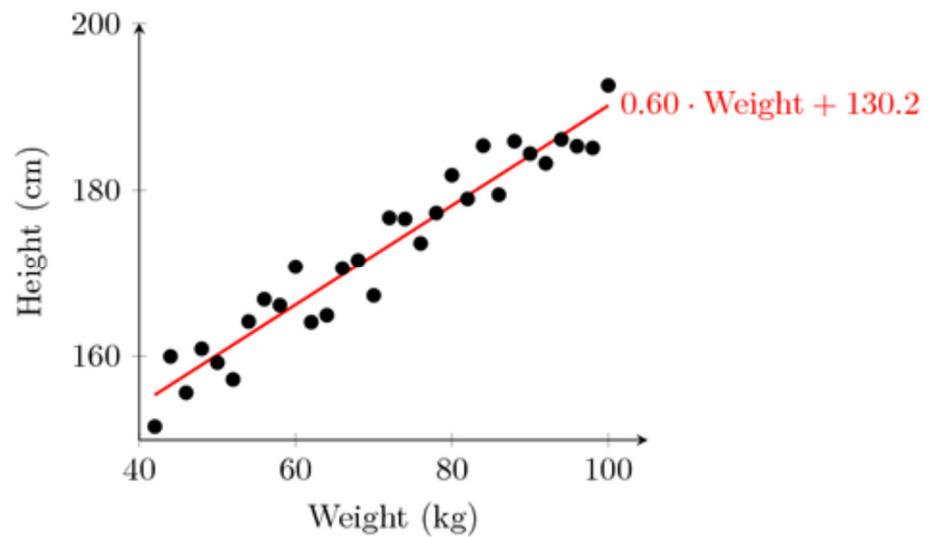
1. Linear Regression
2. Weighted Linear Regression
3. Non-Linear Regression
4. Logistic Regression (Optional)

1. Linear Regression

- ❖ Fit a line to a data set of observations
- ❖ Use this line to predict unobserved values
- ❖ **Goal:** calculate a regression line that's as close to every dot as possible.



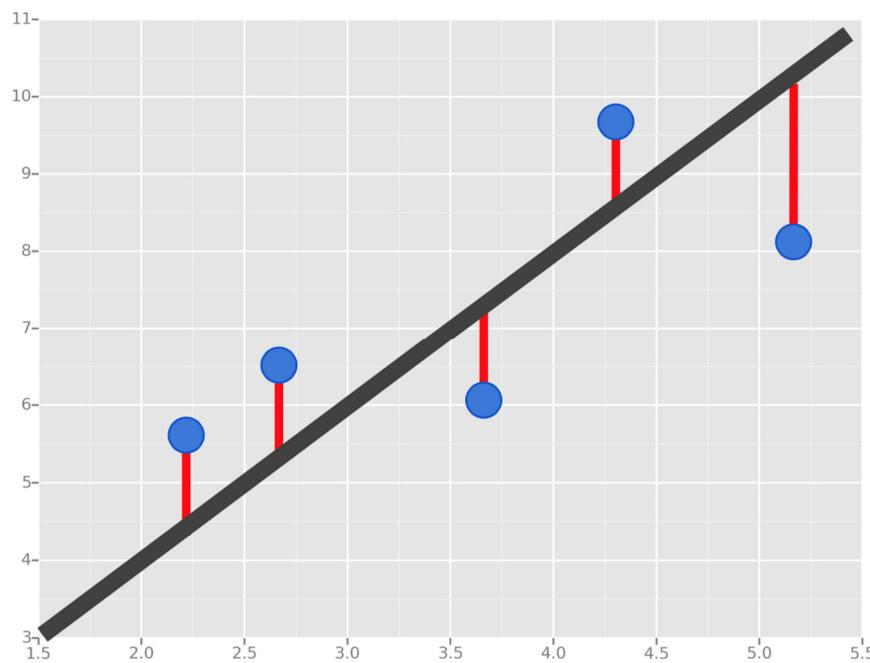
2 points: simple



many data points: not simple

Linear Regression

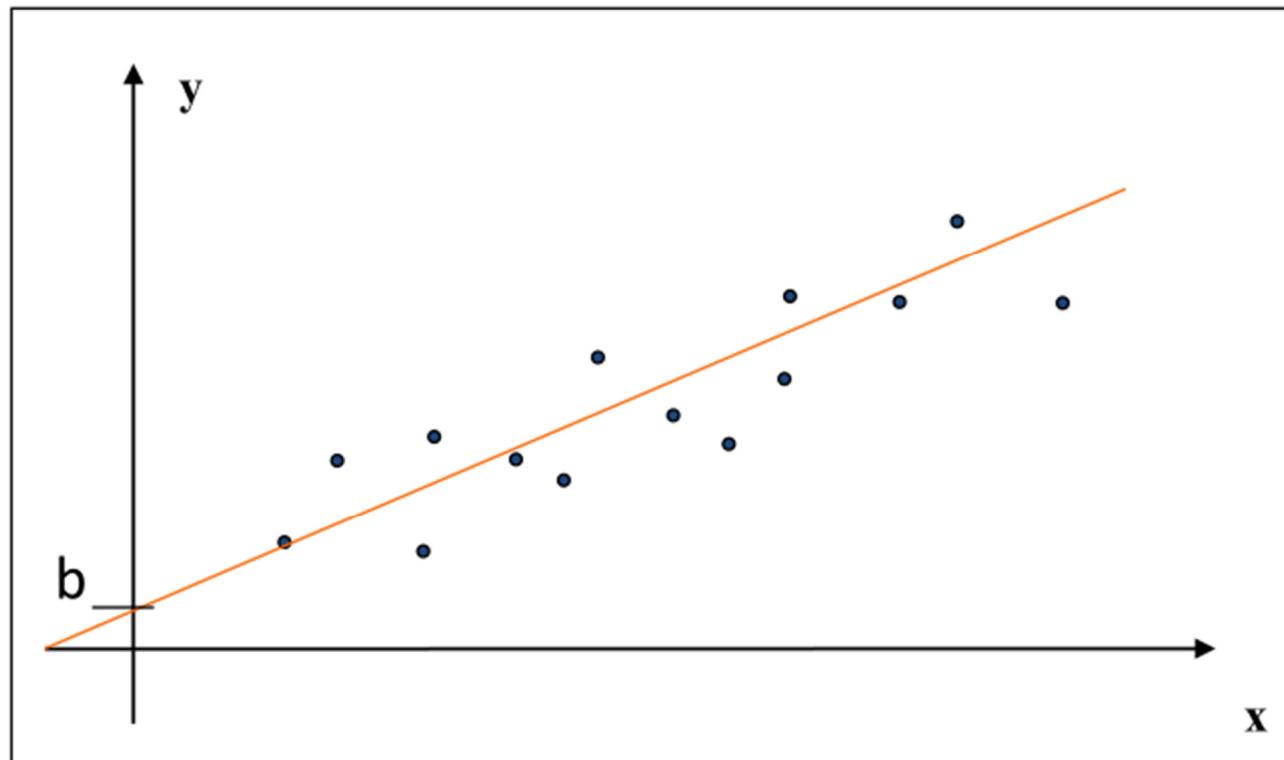
- ❖ **Intuition:** fit a linear regression line that minimizes the vertical distance between all the data points and our line
- ❖ **Distance metrics:**
 - Euclidean distance → aka minimize sum of squares
 - Manhattan distance → aka minimize sum of absolute differences



Linear Regression

- ❖ Linear regression searches a **linear mapping** between input x and output y , parametrized by the slope vector w and intercept b .

$$y = f(x; w, b) = w^T x + b$$

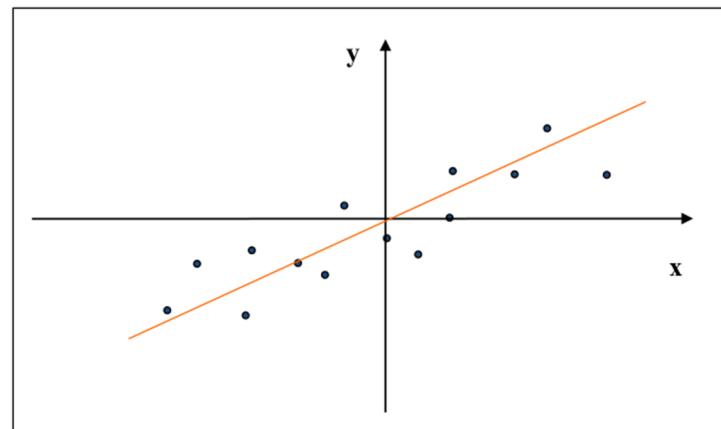


Linear Regression: Optimization

- ❖ **Input:** Pair of M training points $X = [x^1, \dots, x^M]$ and $y = [y^1, \dots, y^M]^T$
 - $x^i \in R^N, y^i \in R$
- ❖ **Optimization Criteria:** least-square error

$$w^* = \min_w \left(\sum_{i=1}^M \frac{1}{2} (w^T x^i + b - y^i)^2 \right)$$

- For convenience, omit the intercept b by **centering** the data:
 - $y' = y - \bar{y}, x' = x - \bar{x}, b' = b + w^T \bar{x} - \bar{y}$
 - ➔ $y' = w^T x'$



For brevity sake, we write $y = w^T x$ with the assumption that the data has been centered

Linear Regression: Optimization

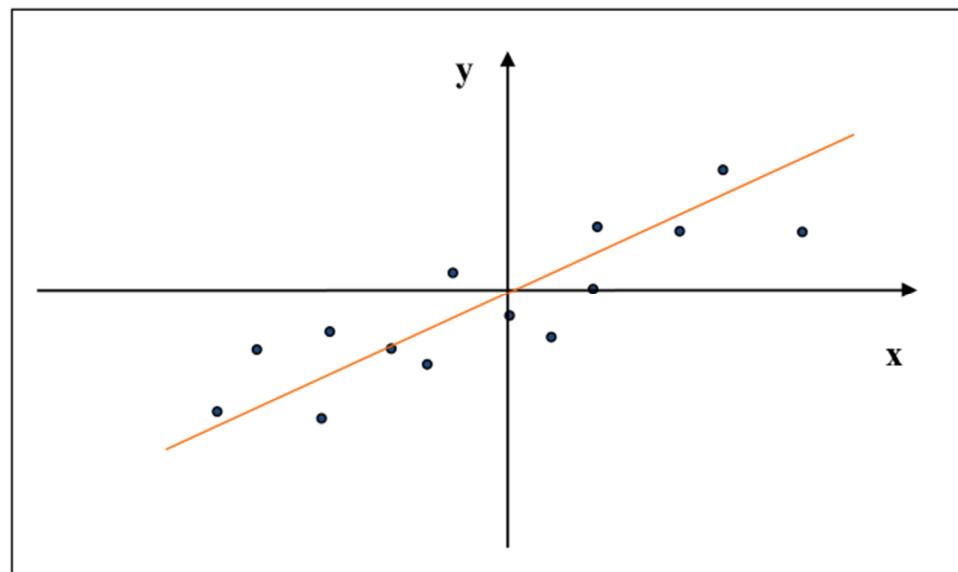
❖ Optimization Criteria: least-square error

- Find the optimal parameter w to **minimize** least-square error:

$$w^* = \min_w \left(\sum_{i=1}^M \frac{1}{2} (w^T x^i - y^i)^2 \right)$$

- This can be solved **analytically** through partial differentiation:

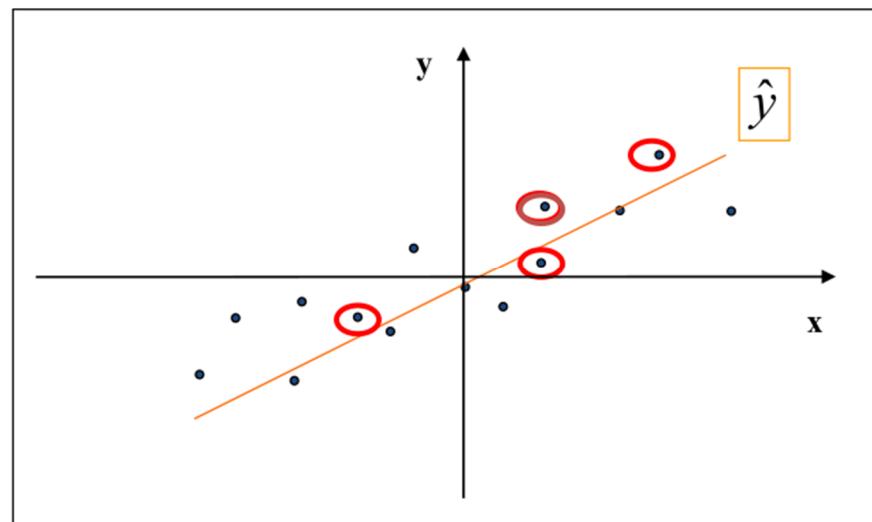
$$w^* = (X X^T)^{-1} X y$$



2. Weighted Linear Regression

- ❖ Sometimes we want to **prioritize** some important points
 - Assume a set of weights β_i for all data points, $\{\beta_1, \dots, \beta_M\}$, $\beta_i \in [0,1]$
 - We set B a diagonal matrix with entries β_i :

$$B = \begin{pmatrix} \beta_1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \beta_M \end{pmatrix}$$



Points in red have large weights.

Weighted Linear Regression

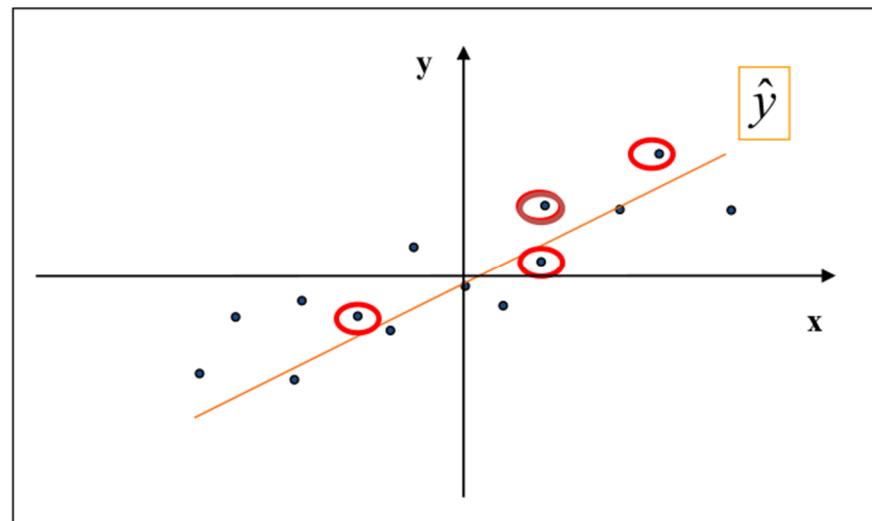
- ❖ **Optimization problem:** weighted least square

$$w^* = \min_w \left(\sum_{i=1}^M \frac{1}{2} \beta_i (w^T x^i - y^i)^2 \right)$$

- ❖ **Solution:**

$$w^* = (Z Z^T)^{-1} Z v$$

➤ where $Z = X B^{1/2}$ and $v = B^{1/2} y$

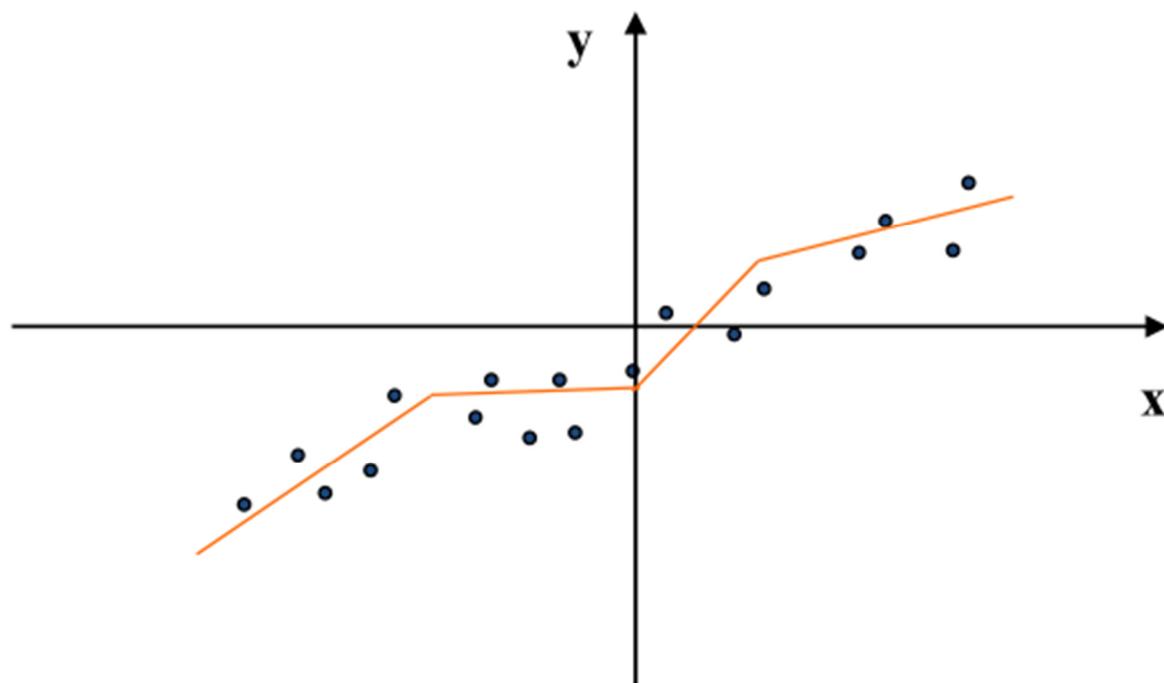


Points in red have large weights.

Locally Weighted Linear Regression

❖ Limitations of (weighted) linear regression

- Set (constant) weights by ourselves is difficult
- Not useful for datasets with **local dependencies** (i.e. different input ranges have different linear dependency)

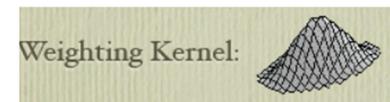
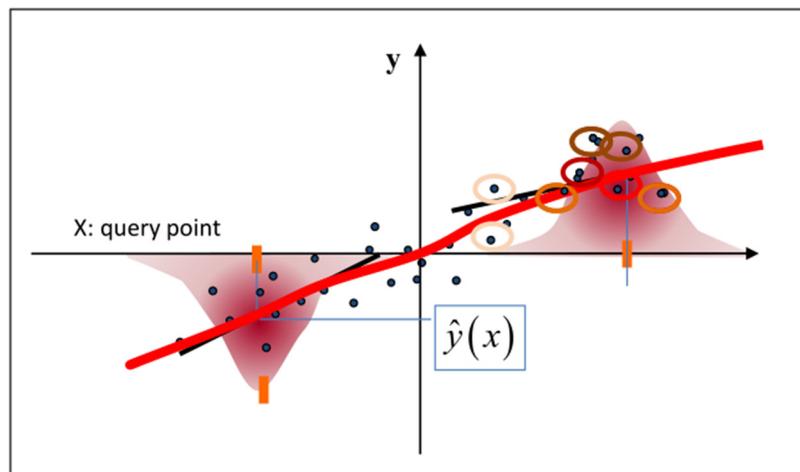


Locally Weighted Linear Regression

- ❖ **New formulation:** estimate is determined through **local influence** of each group of data points

$$\hat{y}(x) = \sum_{i=1}^M \beta_i(x) y^i / \sum_{j=1}^M \beta_j(x)$$

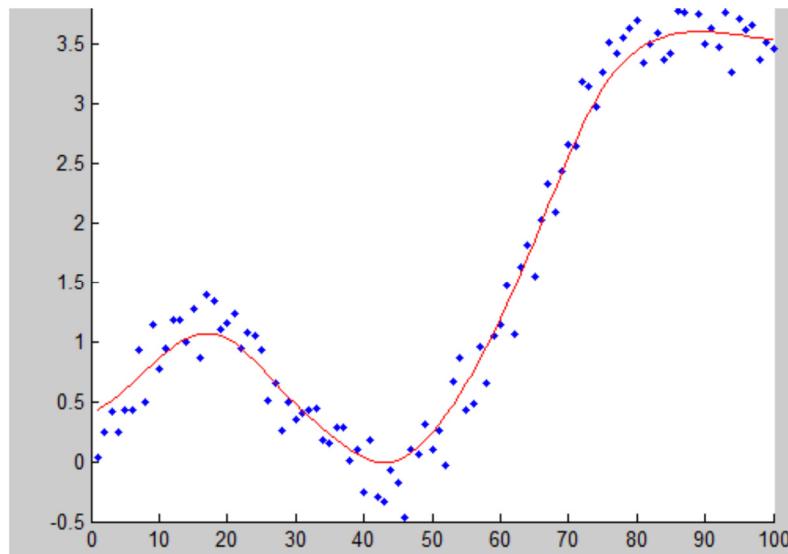
- $\beta_i(x) \in R$: weights function of x
- Now we can generate a smooth function $\hat{y}(x)$ by freely choosing a **weight kernel** $\beta_i(x)$



$$\beta_i(x) = e^{-||x^i - x||}$$

3. Non-linear Regression

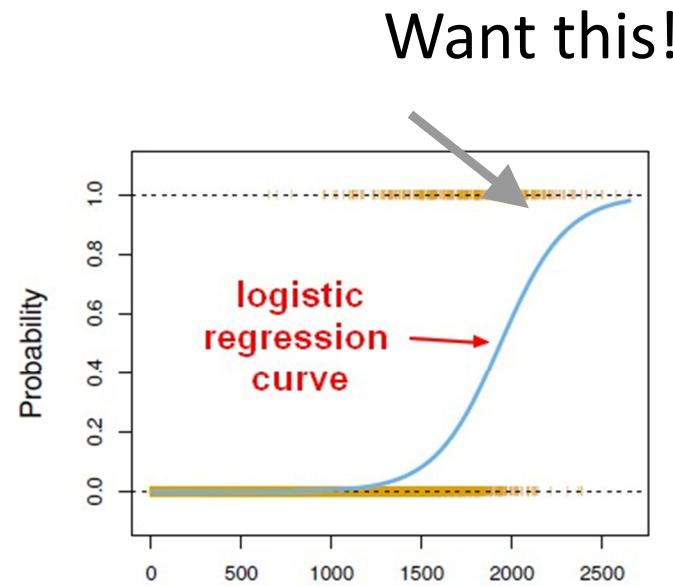
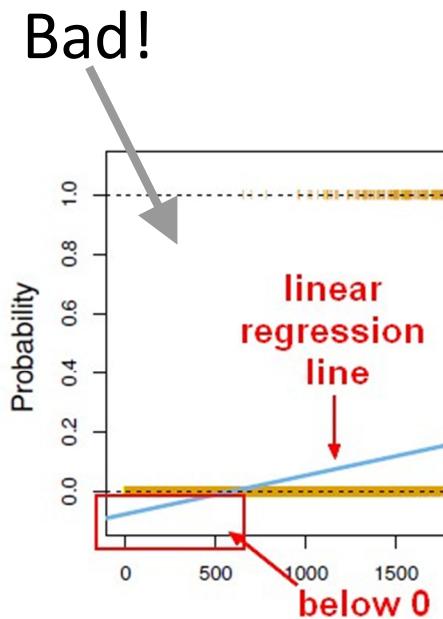
- ❖ Why limit ourselves to straight lines (first order)?
- ❖ Not all relationships are linear
 - E.g. “Second order” polynomial: $y = ax^2 + bx + c$
 - Third order: $y = ax^3 + bx^2 + cx + d$
 - **Higher orders** produce more complex curves



4. Logistic Regression

- ❖ Sometimes we want to see a **degree** of classification instead of yes/no answer
 - Easy for human to investigate false positives and false negatives
- ❖ **Example:**
 - The level of spam of an email
 - The default level of a loan
 - Disease diagnosis
- ❖ **Idea:** linear regression with output is a probability

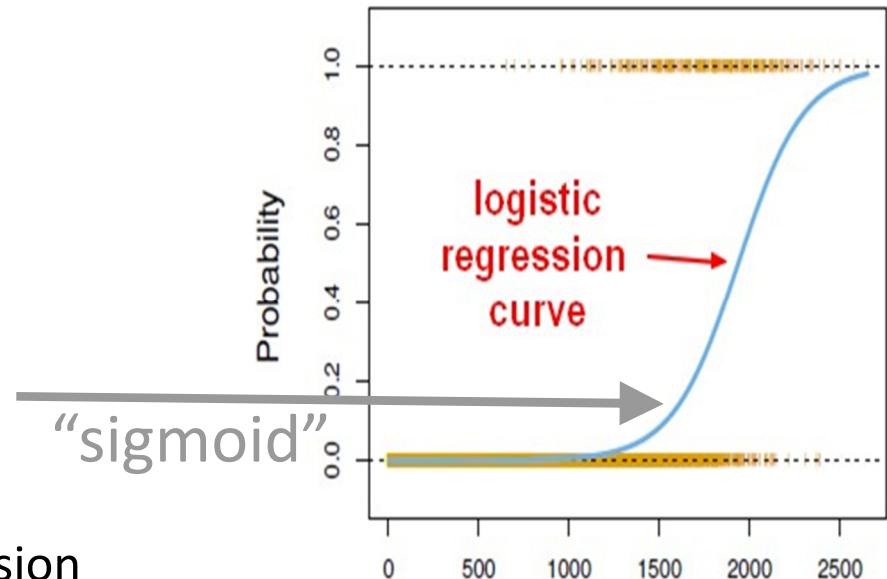
Logistic Regression



- ❖ Trick: don't deal with probabilities, which range from 0 to 1, but with log odds, which range from $-\infty$ to $+\infty$
- ❖ Probability $y \Leftrightarrow$ odds $y/(1-y) \Leftrightarrow$ log odds $\log[y/(1-y)]$
- ❖ Model log odds as a linear function of X

Logistic regression

- ❖ Model log odds as a linear function of X
- ❖ $\beta^T X = \log[y/(1-y)]$
- ❖ Solve for y : $y = 1 / (1 + \exp(-\beta^T X))$
- ❖ Finding best model β :
 - Don't use square loss as in linear regression
 - Use cross-entropy loss instead



Logistic Regression: Pros and Cons

❖ Pros:

- Low variance
- Provide probabilities for outcomes
- Perform similarly to SVM when linear separation

❖ Cons:

- High bias

Multivariate Regression

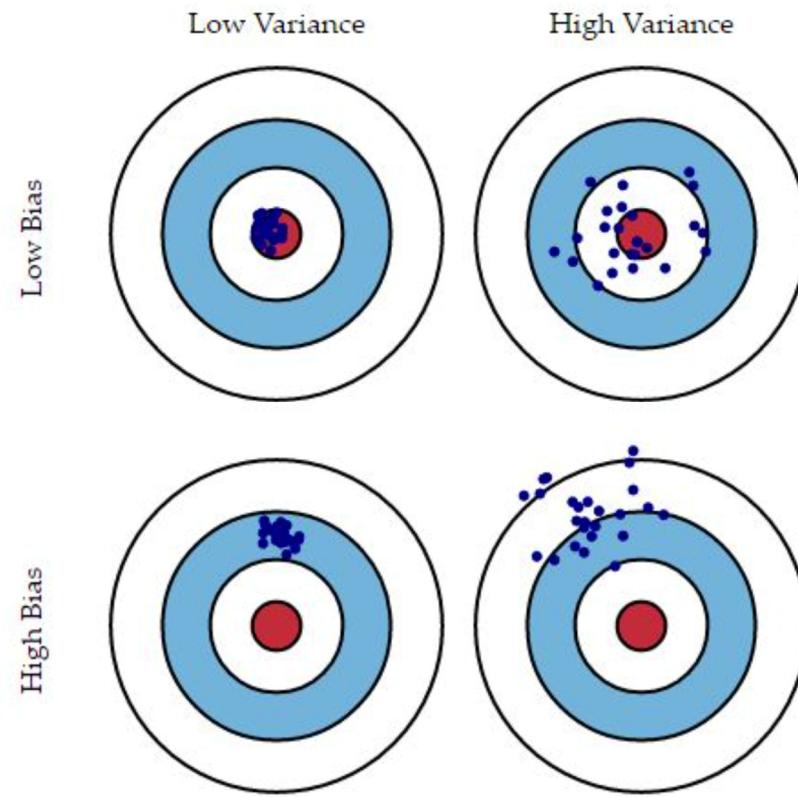
- ❖ What if more than one variable influences the one you're interested in?
- ❖ Example: predicting a price for a car based on its many attributes (body style, brand, mileage, etc.)
- ❖ Model:
 - $\text{Price} = a + b_1 * \text{mileage} + b_2 * \text{age} + b_3 * \text{doors}$
 - Still use least-squares optimization
 - Assumption: different factors are not dependent on each other



IV. Model Evaluation

❖ Bias and Variance:

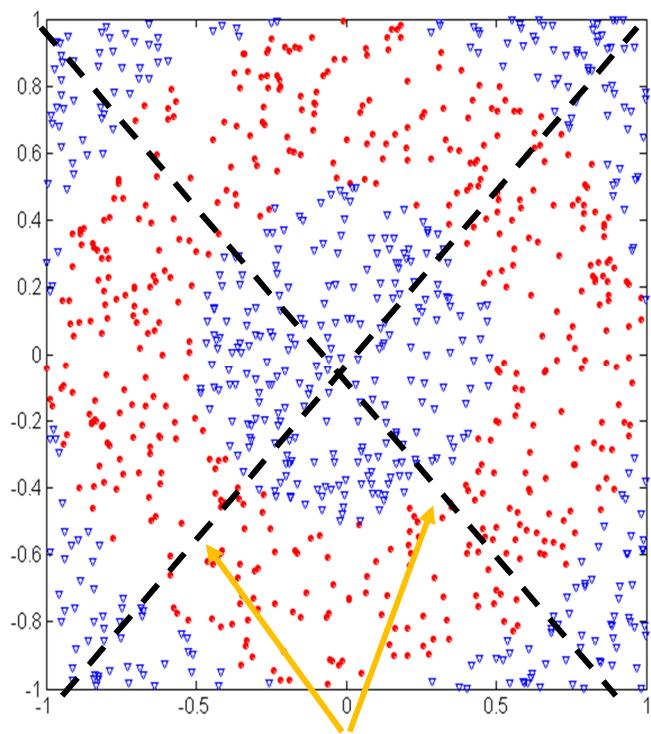
- Bias is how far the mean of your predicted values is from the “real” answer
- Variance is how scattered your predicted values are from the “real” answer



(assuming the center is the correct result)

Underfitting vs. Overfitting

Underfitting

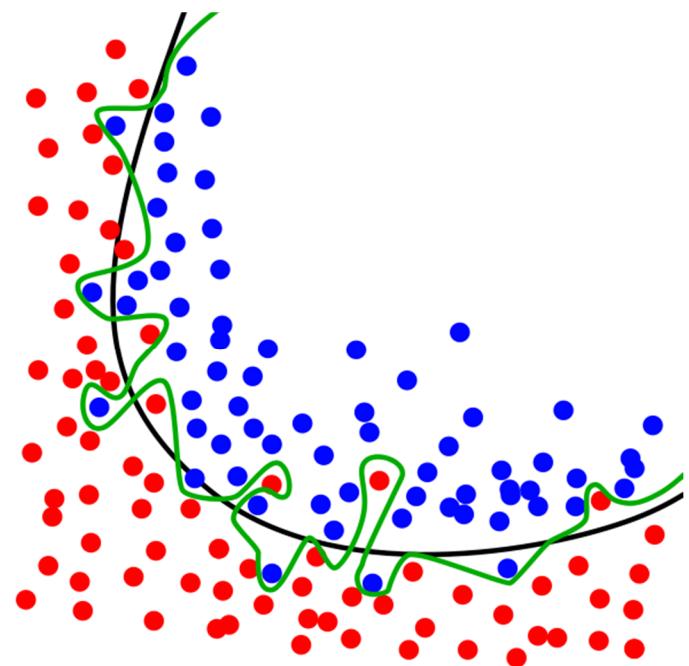


Models are too simple!

Low variance, high bias

E.g. Regression

Overfitting



- Green line is overfitting (tailored too much to the training data)
- Black line is what we want

High variance, low bias

Underfitting vs. Overfitting

❖ Deal with underfitting:

- Increase model complexity (e.g. use polynomial instead of linear)
- Use more training data

❖ Deal with overfitting:

- The more features the better?
 - NO!
 - More features mean less bias, but more variance
- Carefully selected features can improve model accuracy
 - E.g., keep features that correlate with the label y
 - Forward/backward feature selection
 - Regularization (e.g., penalize norm of weight vector)
- K-fold cross validation

K-Fold Cross validation

- ❖ Prevent you from over-fitting a single train/test split
- ❖ General process:
 - Split your data into K randomly-assigned segments
 - Reserve one segment as your test data
 - Train on each of the remaining K-1 segments and measure their performance against the test set
 - Take the average of the K-1 performance measures

Classification: Performance Issues

- ❖ Predictive performance (accuracy, AUC/ROC, precision, recall, F1-score, etc.)
- ❖ Speed and scalability
 - Time to build the model
 - Time to use the model
 - In memory vs. on disk processing
 - Communication cost
- ❖ Robustness
 - Handling noise, outliers, missing values
- ❖ Interpretability
 - Understanding the model and its decisions (black box vs. white box)
- ❖ Compactness of the model
 - Mobile and embedded devices

Classification Performance Measures

- ❖ Classification accuracy can be unreliable when assessing unbalanced data.
 - For example, if there are 95 samples from class A and only 5 from class B in the data set, a particular classifier might classify all the observations as class A.

- ❖ **Confusion matrix** is a table with two rows and two columns that reports the number of false positives (FP), false negatives (FN), true positives (TP), and true negatives (TN). In multi-class classification, a confusion matrix can be computed for each class.

❖ Sensitivity (recall):
$$\frac{TP}{\text{the number of real positives}} = \frac{TP}{TP+FN}$$

❖ Specificity:
$$\frac{TN}{\text{the number of real negatives}} = \frac{TN}{TN+FP}$$

❖ Precision:
$$\frac{TP}{TP+FP}$$

❖ F1-score:
$$\frac{2}{\frac{1}{\text{precision}} + \frac{1}{\text{recall}}} = \frac{2\text{precision} \times \text{recall}}{\text{precision}+\text{recall}}$$

- ❖ Accuracy:

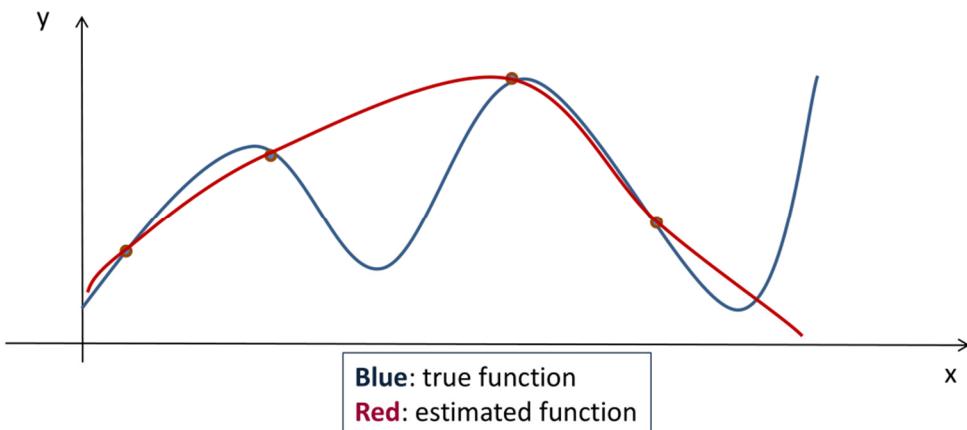
		Actual class		
		Cat	Dog	Rabbit
Predicted class	Cat	5	2	0
	Dog	3	3	2
	Rabbit	0	1	11

confusion matrix for “cat” class

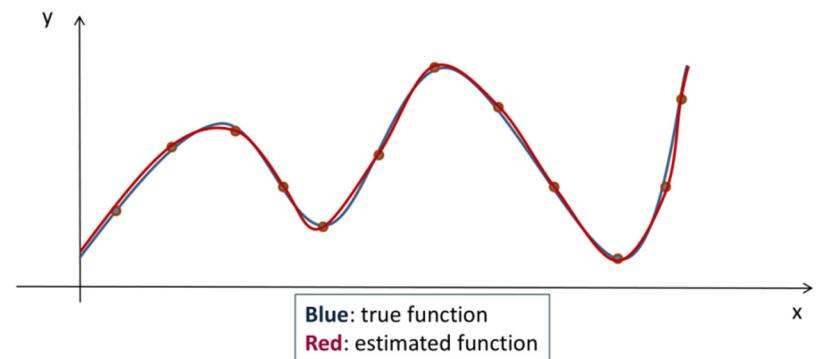
		Actual class	
		Cat	Non-cat
Predicted class	Cat	5 True Positives	2 False Positives
	Non-cat	3 False Negatives	17 True Negatives

Regression: Performance Issues

- ❖ Good prediction depends on the choice of datapoints
- ❖ The more datapoints, the better the fit
- ❖ Computational costs increase dramatically with number of datapoints



Less data points



More data points

References

- [1] <https://blog.exploratory.io/find-correlation-or-similarity-among-categories-or-variables-4813130f53c0>
- [2] <https://blog.exploratory.io/an-introduction-to-regression-analysis-in-exploratory-9422237c0ff8>
- [3] <https://www.wiley.com/en-us/Statistical+Data+Analysis+Explained%3A+Applied+Environmental+Statistics+with+R-p-9780470985816>
- [4] <https://www.bigskyassociates.com/blog/bid/356764/5-Most-Important-Methods-For-Statistical-Data-Analysis>
- [5] <https://www.infogix.com/blog/differentiating-between-cluster-analysis-and-factor-analysis/>
- [6] https://en.wikipedia.org/wiki/Correlation_clustering
- [7] <http://dataaspirant.com/2014/09/27/classification-and-prediction/>
- [8] <http://syllabus.cs.manchester.ac.uk/ugt/2017/COMP24111/lectures.php>
- [9] <http://www.r2d3.us/visual-intro-to-machine-learning-part-1/>
- [10] https://github.com/ctufts/Cheat_Sheets/wiki/Classification-Model-Pros-and-Cons