

CS5014: Machine Learning, Revision Notes

Gergely Flamich

May 12, 2018

1 Overview

In ML, we usually have some data in the following general setup:

- **Features:** literally any sort of observations from the real world, collected conveniently into a matrix:

$$X = \begin{bmatrix} x_{1,1} & x_{1,2} & \dots & x_{1,n} \\ x_{2,1} & x_{2,2} & \dots & x_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ x_{m,1} & x_{m,2} & \dots & x_{m,n} \end{bmatrix}$$

where each column $\mathbf{x}^{(i)}$ contains m observations of the i -th feature

- **Labels:** In supervised learning we also have another set of data, called the ground truth of the object we are trying to predict from the features. These are conveniently collected into a vector \mathbf{y} with m elements, each y_i corresponding to a set of observations, row \mathbf{x}_i in X .

The general assumption for supervised learning is that each feature is associated with a random variable X_i over some distributions and the ground truth is associated with a random variable Y over some distribution. It is usually assumed that we know the distributions of X_i and Y is hidden. Finally, we assume that there is some function f such that

$$Y = f(\mathbf{X}, \theta) + \epsilon$$

for some model parameters θ and some small error term ϵ . Then, the way we will perform prediction is by using the distribution

$$\hat{Y} = f(\mathbf{X}, \theta).$$

f is usually referred to as the **model function**.

1.1 Machine Learning Pipeline

1. Obtain data
2. Clean data
3. Visualise & Analyse data
4. Prepare data: augmentation, feature selection etc.
5. Perform the learning task

6. Evaluate algorithm

7. Test algorithm

1.2 Flavours of ML

Supervised learning:

- Output is continuous: regression
- Output is discrete: classification

Unsupervised learning: Clustering

1.3 Common issues in ML

- **Overfitting - High Variance:** The model is too specific to the data. It fits it very well, but doesn't generalise well.
- **Underfitting - High Bias:** The model is too general, it doesn't explain the training data well, and it doesn't explain new data that well either.
- **Data Leakage:** Some information from the test set "leaks" into the training/ validation set, and so the model will be biased to the seen data.
- **Dataset bias:** The training or validation set is not representative of the general domain of the problem.

2 Data Preparation Methods

- **Cleaning:** Remove nonsensical/null values, put data in the right format
- **Feature Extraction:** use raw features and combine/ transform them to purify some information to aid learning: requires domain knowledge usually.
- **Feature Selection:** Improve interpretability and computational feasibility. Remove unnecessary features (based on domain knowledge or statistical analysis): univariate tests, forward stepwise selection, backward stepwise selection etc.
- **One-Hot Encoding:** For classification tasks with k classes, it is commonplace to represent different classes as pairwise orthogonal unit vectors in \mathbb{R}^k . This ensures that $\|c_i - c_j\| = \text{constant}$ for every class, i.e. there is no bias towards any class. Example:

$$\begin{bmatrix} c_1 \\ c_4 \\ c_3 \\ \vdots \\ c_2 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & 1 & 0 & 0 \end{bmatrix}$$

- **Decorrelation:** remove linear dependencies from the features: whitening, PCA.

- **Normalisation:**

$$x_i \rightarrow \frac{x_i - \min \mathbf{x}}{\max \mathbf{x} - \min \mathbf{x}}.$$

- **Standardisation:** if X_k is assumed to be symmetrically, or even better, normally distributed with mean μ_k and STD σ_k , this gives better results than normalisation:

$$x_i \rightarrow \frac{x_i - \mu_k}{\sigma_k}$$

- **Data Augmentation:** artificially generate more data from our preexisting set by introducing small modifications to it.

3 General ML Methodology

3.1 Regularisation

Purpose: Improve the generalisability of the model and prevent overfitting.

General Idea: Penalise an individual parameter growing too large.

Example: Squared Loss

$$L(\theta) = \|X\theta - \mathbf{y}\|^2 \rightarrow J(\theta) = \|X\theta - \mathbf{y}\|^2 + \lambda \|\theta\|^2$$

for some **regularisation parameter** $\lambda \in \mathbb{R}^+$. If we regularise using the L_2 norm of θ , then the algorithm is called **ridge regression**, if we use the L_1 norm, then it's called **lasso**.

3.2 Model Evaluation

Train-Validation-Test Split Common ratios: 60 - 20 - 20, 80 - 10 - 10, etc.

- **Training data:** we train our model using the training data
- **Validation data:** we evaluate a trained model on the (preferably) unseen validation data according to some metric(s). Based the model performance, we can tune hyperparameters/choose a different model.
- **Test data:** Should always be unseen data to avoid bias in the trained model towards previously seen data. If this is not done, it constitutes data leakage and our results will be meaningless. The results on the test data are the results that we usually report, because it tells us how well our model generalises.

Cross Validation : If we don't have enough data to perform a proper train-validation-test split, we perform cross-validation. We perform a train-test separation, say 80-20. Then, taking the training data, we split into k equal portions. Then we train our models k times on $k - 1$ different portions and validate on the one that was left out. Then, this helps mitigate some of the bias in our dataset. We call this k-fold CV.

Stratified Sampling : ensures consistent representation of classes in training and validation sets, by sampling each subpopulation (class) of the dataset independently, and then merging the results.

Performance Metrics : Regression:

- Squared error:

$$SE = \|\mathbf{y} - \hat{\mathbf{y}}\|_2^2$$

- Absolute error:

$$AE = \|\mathbf{u} - \hat{\mathbf{y}}\|_1$$

- Coefficient of determination: (sometimes called explained variance)

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

- Root Mean Squared Error:

$$RMSE = \sqrt{\frac{SE}{m}}$$

RMSE is useful, because it has the same units as the data, so we can directly observe how well our model is performing.

(Binary) Classification:

- Accuracy:

$$\frac{TP + TN}{P + N}$$

- Error rates - False Positive Rate:

$$FPR = \frac{FP}{TN + FP}$$

False Negative Rate:

$$FNR = \frac{FN}{FN + TP}$$

- Specificity:

$$1 - FPR$$

- Sensitivity:

$$1 - FNR$$

- Precision:

$$\frac{TP}{TP + FP}$$

- Recall:

$$\frac{TP}{TP + FN}$$

- F1 - score:

$$2 \frac{Prec \cdot Recall}{Prec + Recall}$$

The above pairings of classification metrics are “working against each other”: one of them increasing will mean a decrease of their pair.

ROC Curve: Receiver Operation Characteristic Curve: TPR vs FPR

Precision-Recall Curve: duh

Average Precision (AP): integral of the precision recall curve

Metrics for multiclass classification:

- Report One-vs-All scores
- Confusion Matrix: for k classes it is a $k \times k$ matrix $C = [c_{ij}]$, where $c_{ij} = \#$ of how many times an object of class c_j was predicted to be class c_i .

4 Regression Methods

4.1 Linear Regression

One of the simplest models f relies on the further assumption that Y depends linearly on \mathbf{X} . In particular,

$$\hat{Y} = f(\mathbf{X}, \theta) = \mathbf{X}\theta + b$$

where b is called the **bias**. No Then, we can quantify the goodness-of-fit of the model, by the squared error, specified by

$$L(\theta, b) = \|\mathbf{y} - f(X, \theta, b)\|^2$$

Then, to find the optimal set of parameters theta, we can differentiate

$$\begin{aligned}\frac{\partial L}{\partial \theta} &= \frac{\partial}{\partial \theta} \|\mathbf{y} - f(X, \theta, b)\|^2 \\ &= 2 \|\mathbf{y} - f(X, \theta, b)\|^2 \frac{\partial}{\partial \theta} \|\mathbf{y} - f(X, \theta, b)\| \\ &= 2 \|\mathbf{y} - f(X, \theta, b)\| \frac{(\mathbf{y} - f(X, \theta, b))^T}{2 \|\mathbf{y} - f(X, \theta, b)\|} \frac{\partial}{\partial \theta} \mathbf{y} - f(X, \theta, b) \\ &= (\mathbf{y} - X\theta)^T \cdot -X\end{aligned}$$

Now to find the minimum, set $\frac{\partial L}{\partial \theta} = 0$.

$$\begin{aligned}0 &= -\mathbf{y}^T X + \theta^T X^T X \\ \mathbf{y}^T X &= \theta^T X^T X \\ \mathbf{y}^T X (X^T X)^{-1} &= \theta^T\end{aligned}$$

Hence

$$\theta = (X^T X)^{-1} X^T \mathbf{y}$$

The above is called the **normal equation**. An alternative way to solve the above problem is **gradient descent**. Iterate

$$\theta \leftarrow \theta - \alpha \cdot \nabla L$$

for some appropriate $\alpha \in \mathbb{R}$. Then it can be shown that the θ converges to the minimal solution. α is called the **learning rate**.

Regularised Normal Equation

$$\theta = (X^T X + \lambda I)^{-1} X^T$$

4.1.1 Basis Expansion

Provided we have k real features, we can modify our space in which we consider the linearity of our model. We can transform a set of observations \mathbf{x} into the new space through an arbitrary function $h : \mathbb{R}^k \rightarrow \mathbb{R}^n$ by calculating $h(\mathbf{x})$. then, our model changes to:

$$f_\theta(X) = X\theta \quad \rightarrow \quad f_\theta(X) = h(X)\theta.$$

Crucially, f_θ is still linear in the space defined by h , but the components of h may be non-linear!

Special Case: Polynomial regression: the h_i components of h are different combinations of the input features.

4.2 Logistic Regression

Model the classification problem using the **logistic function**:

$$\sigma_\theta(\mathbf{x}) = \frac{1}{1 + e^{-\theta^T \mathbf{x}}} \quad \text{where } \mathbf{x}, \theta \in \mathbb{R}^{n \times 1}$$

It is easy to see that $\sigma_\theta(\mathbf{x}) \geq 0.5$ iff $\theta^T \mathbf{x} \geq 0$ and less than 0.5 otherwise.

$$B = \{\mathbf{x} \mid \sigma_\theta(\mathbf{x}) = 0.5\}$$

is called the **decision boundary** of the model.

This model allows for **binary classification**

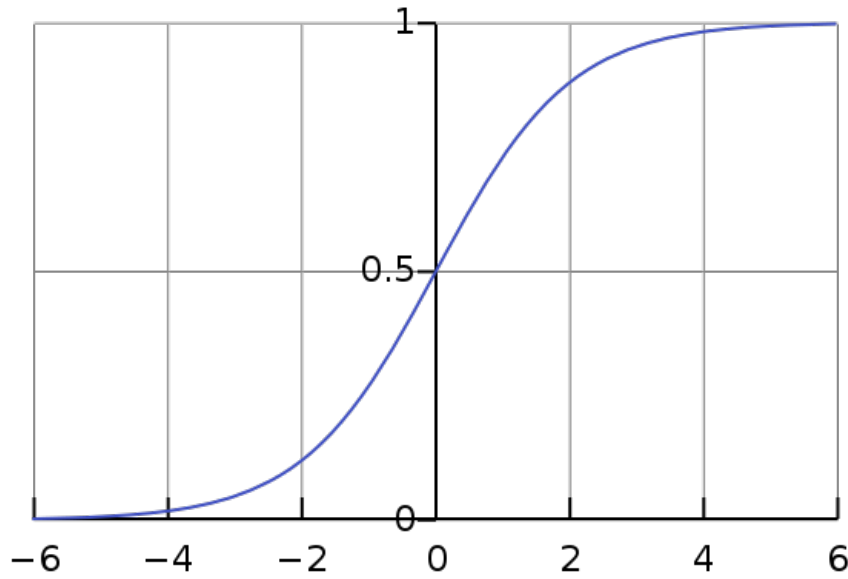


Figure 1: Logistic function $\sigma_\theta(\mathbf{x})$

The loss function for logistic regression is rooted in information gain, so L is defined to be the cross-entropy of the data:

$$L(\theta) = - \sum_{i=1}^m y_i \log(f_\theta(\mathbf{x}_i)) + (1 - y_i) \log(1 - f_\theta(\mathbf{x}_i))$$

Optimise this with gradient descent as before.

4.2.1 Multi-Class Logistic Regression

One-vs-All : For k classes train k classifiers: for each class c_i one class is c_i and one class is everything else. Then, to predict, take prediction with highest confidence.

One-vs-One : For k classes train $\binom{k}{2}$ classifier: one for each pair of c_i, c_j classes. Then, take a majority vote.

5 Bayesian Decision Theory

Main idea of Bayesian Decision Theory: minimise expected loss.

5.1 MLE

Assume set of distributions: $\{p_\theta \mid \theta \in \Theta\}$.

Assume the data is an i.i.d. sample of the above family: $X_1, \dots, X_n \sim p_\theta$.

Then, attempt to estimate the true θ that generated the sample:

$$\theta_{MLE} = \arg \max_{\theta \in \Theta} \mathbb{P}[X \mid \theta]$$

where $\mathbb{P}[X \mid \theta]$ is called the **likelihood function**.

Pros

- Easy to compute
- Interpretable
- Nice asymptotic properties
- Invariant under reparametrisation:

$$g(\theta_{MLE}) = (g(\theta))_{MLE}$$

Cons

- Point estimate
- Overfitting

5.2 MAP

Similar assumptions to MLE. Instead of the likelihood function, estimate the **posterior distribution** of θ . We wanna estimate $\mathbb{P}[X \wedge \theta]$, we can use Bayes' Rule:

$$\mathbb{P}[\theta \mid X] = \frac{\mathbb{P}[X \mid \theta] \mathbb{P}[\theta]}{\mathbb{P}[X]}$$

here $\mathbb{P}[X]$ is constant, so we need to minimise

$$\theta_{MAP} = \arg \max_{\theta \in \Theta} \mathbb{P}[\theta \mid X] = \arg \max_{\theta \in \Theta} \mathbb{P}[X \mid \theta] \mathbb{P}[\theta]$$

Pros

- Easy to compute
- Interpretable
- Avoids overfitting by putting a prior on θ - Regularisation
- Converges to MLE

Cons

- Point estimate
- Not invariant under reparametrisation:

$$g(\theta_{MAP}) \neq (g(\theta))_{MAP}$$

- Must assume prior on θ

5.3 Density Estimation

Local density estimation assume that the density of the data is locally constant. Then for some volume V around some point $\mathbf{X} = [X_1, \dots, X_P]^T$, if we have K samples in the volume out of a total of N samples, then define

$$\mathbb{P}[X] = \frac{K}{NV}.$$

How we pick K and V determines the method.

Histograms: split up the space into a grid of cubes called **bins**. Then, we fixed V and we'll get K by counting up the samples in every bin.

Pro: Cheap

Con: often too crude.

Parzen density estimation: Given some point X , assume that the density is some local function of nearby samples. So define:

$$\mathbb{P}[X] = \frac{1}{N} \sum_{n=1}^N \frac{1}{V} \cdot K\left(\frac{X - x_n}{h}\right)$$

where h is the window size, i.e. it controls the locality of the **kernel function** K . The kernel functions must be symmetric probability density functions themselves.

Examples:

- Tophat kernel: K is a uniform density function between $-h/2$ and $h/2$.
- Gaussian: $K = \mathcal{N}(0, h/2)$.
- Exponential
- Linear

5.4 K-Nearest Neighbours

If we fix K and let V grow according to

$$\mathbb{P}[X] = \frac{K}{NV}$$

then we can estimate class likelihoods according to Bayes' Rule:

$$\mathbb{P}[C_i | X] = \frac{\mathbb{P}[X | C_i] \mathbb{P}[C_i]}{\mathbb{P}[X]} = \frac{K_i}{K}$$

6 Maximal Margin Classifiers & SVMs

Main idea: Instead of best fit to all of data, try to **separate** it as best as possible, so focus more on data close to the margin. So, maximise margin separating the data.

6.1 Setup

Decision function:

$$\hat{y} = \text{sgn}(\mathbf{x}^T \beta + \beta_0)$$

Distance to the margin:

$$y_i(\mathbf{x}^T \beta + \beta_0) \geq 0 \quad \forall i \in \{1, \dots, N\}$$

Then, maximise this distance by enforcing it to be greater than some margin M :

$$\frac{1}{\|\beta\|} y_i(\mathbf{x}^T \beta + \beta_0) \geq M$$

A **support vector** is any vector that is directly on the margin, i.e. the above inequality is an equality. Setting $M = \frac{1}{\|\beta\|}$:

$$y_i(\mathbf{x}^T \beta + \beta_0) = 1$$

This is a **hard margin**, because this is only possible to be solved if the data is **separable** and it fails if it isn't.

6.2 Soft Margin Relaxation

Allow a few points be on the wrong side. Introduce some slack variables $\xi = [\xi_1, \dots, \xi_N]$, that will allow us to let data points cross the margin:

$$y_i(x_i^T \beta + \beta_0) \leq 1 \quad \rightarrow \quad y_i(x_i^T \beta + \beta_0) \leq 1 - \xi_i$$

And then constrain ξ to not grow too big:

$$\xi_i \geq 0, \quad \|\xi\|_1 \leq \text{const}$$

Note: $\xi_i \geq 0$: \mathbf{x}_i is inside margin. $\xi_i > 1$: \mathbf{x}_i is misclassified. The above can be reformulated to a cost function as:

$$J(\beta) = \frac{1}{N} \sum_{i=1}^N \underbrace{(1 - y_i(x_i^T \beta + \beta_0))_+}_{\text{hinge loss}} + \underbrace{\lambda \|\beta\|_2^2}_{\text{regularisation}}$$

Note: The correctly classified data points don't contribute to the loss at all.

Note: Hinge loss is not differentiable at the hinge point: use **subgradient descent**: pick subgradient that minimises loss the most.

6.3 Kernel Trick

Turn the soft margin classification problem into an objective function through **Lagrange relaxation**:

$$L_p = \underbrace{\frac{1}{2} \|\beta\|^2}_{\text{optimisation objective}} + \underbrace{C \sum_{i=1}^N \xi_i}_{\text{slack variable constraint}} - \underbrace{\sum_{i=1}^N \alpha_i (y_i(\mathbf{x}_i^T \beta + \beta_0) - (1 - \xi_i))}_{\text{maximal margin condition}} - \underbrace{\sum_{i=1}^N \mu_i \xi_i}_{\text{non-negativity condition}}$$

where C, α_i, μ_i are called **Lagrange multipliers**, they represent the hard constraints. This is called the **primal problem** and hard to optimise: turn it into the **dual problem**. We differentiate with respect to β, β_0, ξ_i and set the derivatives to 0 since we know that any solution must be a minimum so

$$\beta = \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i \quad 0 = \sum_{i=1}^N \alpha_i y_i \quad \alpha_i = C - \mu_i$$

Using these constraints and putting it back into L_P and simplifying, we get the **dual function**:

$$L_D = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \underbrace{x_i^t x_j}_{\text{dot product}}$$

We can use the dual function to optimise β . Similarly, classification, we have

$$\beta_0 = \sum_{i=1}^N \hat{\alpha}_i y_i \langle \mathbf{x}, \mathbf{x}_i \rangle$$

We can perform learning and classification if we know the dot product. The idea then, is to replace the original dot product with different ones: these are called **kernels**. Examples:

- Polynomial kernel: corresponds to basis expansion

$$K(\mathbf{x}, \mathbf{x}') = (1 - \langle \mathbf{x}, \mathbf{x}' \rangle)^d$$

- Radial Basis Function: Considers the distance as distances of normal distributions

$$K(\mathbf{x}, \mathbf{x}') = \exp(-\gamma \|\mathbf{x} - \mathbf{x}'\|^2)$$

- Neural Network Kernel

7 Neural Networks and Deep Learning

Feature selection is difficult and requires a lot of domain knowledge: how about we let the machine do it?

- Have $L \geq 2$ layers: first layer is the **input layer**, last layer is the **output layer**. Everything in between is a **hidden layer**.
- Each layer has some $\mathbf{z}_i \in \mathbb{R}^{m_i}$ units. Then, each unit on some layer is connected to every unit from the previous layer in a linear fashion:

$$\forall 0 \leq i \leq L-1 : \quad \Theta^{(i)} \in \mathbb{R}^{m_{i-1} \times m_i} : \quad \mathbf{z}_i = \mathbf{a}_{i-1}^T \Theta^{(i)} + \mathbf{b}_i.$$

where \mathbf{b}_i are called the bias of the layer. Then, we calculate the **activations** of the units:

$$\mathbf{a}_i = g(\mathbf{z}_i)$$

for some $g : \mathbb{R} \rightarrow \mathbb{R}$. Usual picks for g :

- Logistic function: σ

- Hyperbolic tangent: \tanh
- Rectified Linear Unit (ReLU): $g(x) = x_+$

- Then, a single **forward propagation** of the network can be written as:

$$\hat{\mathbf{y}} = g(\dots g(g(\mathbf{x}^T \Theta^{(1)})^T \Theta^{(2)})^T \dots \Theta^{(L)})$$

This can be used to perform regression or classification.

- The cost function is

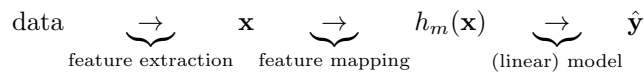
$$J(\Theta) = \frac{1}{m} \sum_{i=1}^m L(\mathbf{y}, \hat{\mathbf{y}}) + \frac{\lambda}{2m} \sum_{l=1}^L \|\Theta^l\|^2$$

where the loss function L can be cross-entropy loss, squared error, etc.

- Then, to learn, we minimise J using gradient descent. To calculate the derivatives of $\Theta_{i,j}^{(l)}$, we need to calculate them for the layer after them, and then use the chain rule.
- It is crucial to randomly initialise the weights of the network to prevent two units learning the same thing in the network.
- During implementation, use gradient checking (numerically approximate the derivatives and compare them with analytical results).

7.1 Deep Learning

Classical ML pipeline:



Deep Learning: Any algorithm that attempts to learn intermediate representations of the data.

Problems of DL:

- Computational Complexity
- Availability of data: we need a lot of it to learn a lot of parameters correctly
- Overfitting
- (Highly) Non-convex cost function
- Vanishing gradient problem: if we have a sigmoid function σ , then if $0 \ll |x|$ then $\sigma'(x) \approx \sigma'(x \pm h)$ for small $h > 0$, meaning that it is hard for the network to learn.
Solution: use ReLU, leaky ReLU or Softplus.
- Picking a network topology

7.2 Regularisation techniques

- Get a lot of data: this has the best regularisation effect
- Data augmentation: modify preexisting data sensibly to improve learning outcome: e.g. add noise, rotate image, etc.
- Weight regularisation:

$$J(\Theta) = L(\Theta) + \underbrace{\frac{\lambda}{2} \sum_{\substack{i,j,l \\ i \neq 0}} (\Theta_{i,j}^{(l)})^2}_{\text{regularisation term}}$$

- Early stopping: if the loss function stagnates, stop the learning process to stop overfitting because the **validation loss might increase!**
- Dropout: for each unit, do not consider it with some probability $0 < p < 1$. This adds natural noise, as we're training only on partial data and forces the network to be much more general.
- Stochastic / Minibatch Gradient Descent: instead of calculating gradients based on the entire dataset, introduce noise by shuffling the dataset, splitting it up into smaller batches and calculating gradients based on the smaller batches

Note: we need to decrease the learning rate otherwise we will never converge. This is called **weight decay**.

Note: SGD improves behaviour around saddle points / local minima.

- Batch normalisation: after the activation of every layer, normalise the activations! Speeds up training and has a regularisation effect.

Weight sharing: instead of learning a bunch of different stuff, we could share some weights along nodes and help them optimise the features extracted by those shared weights. Leads to CNNs.

Max-pooling: reduce the dimensionality of the data by throwing away locally insignificant data.

Recent DL successes: AlexNet, ResNet

8 Unsupervised Learning

8.1 K-means clustering

1. Start with K randomly initialised centroids μ_i
2. During each iteration, find the closest centroid to each data point

$$c_i \leftarrow \min_{k \in \{1, \dots, K\}} \|\mu_i - \mathbf{x}_i\|$$

3. set the centroids to the means of their assigned datapoints

$$\mu_k \leftarrow \frac{1}{(\# \text{ of datapoints in cluster } i)} \sum_{c_i=k} x_i$$

4. Repeat 2. and 3. until convergence

To evaluate the clustering algorithm, we can use the loss function

$$J(\mathbf{c}, \mu) = \frac{1}{m} \sum_{i=1}^m \|\mathbf{x}_i - \mu_{c_i}\|^2,$$

i.e. we want to minimise the total mean distance. Then, to find the best clustering, run the algorithm 100 times, say, and use the one that minimises J .

8.2 PCA

Dimensionality reduction: calculate a new orthogonal basis for the data and throw away axes along which the variance is low