

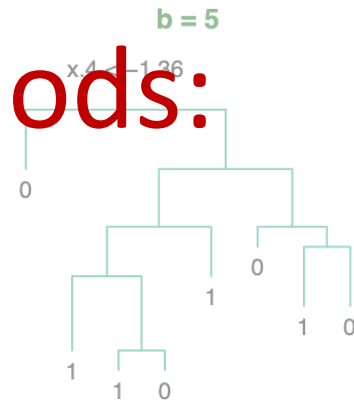
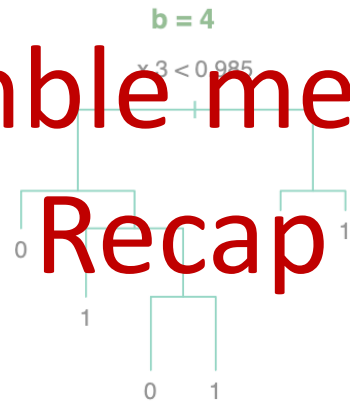
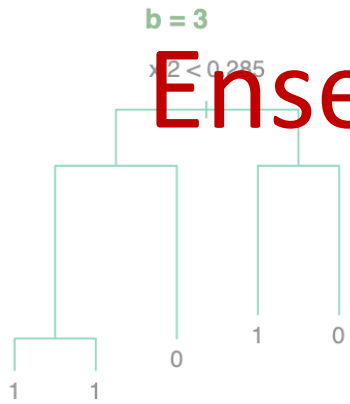
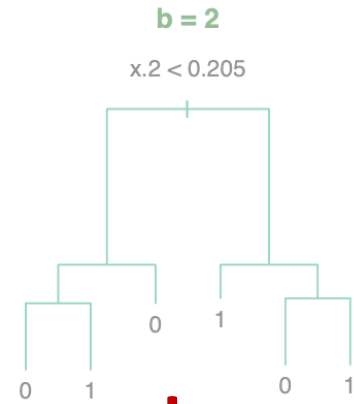
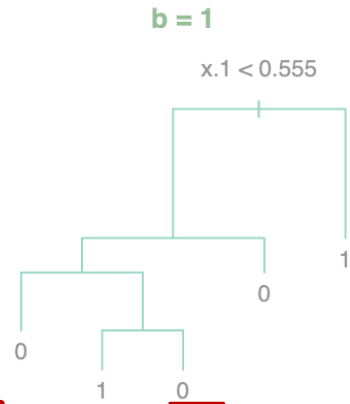
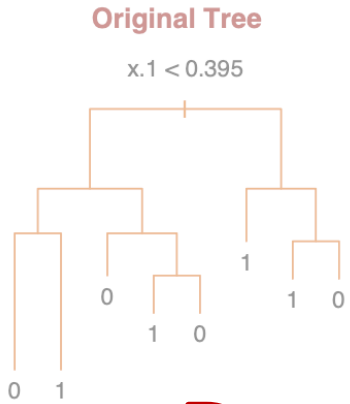
CSCI 567: Machine Learning

Vatsal Sharan
Spring 2024

Lecture 10, Mar 29

Administrivia

- HW4 is due in about two weeks (April 11 at noon).
- Mid-term project check-ins in week of April 15
- Today's plan:
 - Finish ensemble methods
 - Unsupervised learning:
 - PCA



Decision Trees and Ensemble methods:

Recap

Ensemble methods

- Bagging
- Random forests
- Boosting: Basics
- Adaboost
- Gradient boosting

Bagging

Collect T subsets each of some fixed size (say m) by sampling with replacement from training data.

Let $f_t(\mathbf{x})$ be the classifier (such as a decision tree) obtained by training on the subset $t \in \{1, \dots, T\}$. Then the aggregated classifier $f_{agg}(\mathbf{x})$ is given by:

$$f_{agg}(\mathbf{x}) = \begin{cases} \frac{1}{T} \sum_{t=1}^T f_t(\mathbf{x}) & \text{for regression,} \\ \text{sign} \left(\frac{1}{T} \sum_{t=1}^T f_t(\mathbf{x}) \right) = \text{Majority Vote} \{ f_t(\mathbf{x}) \}_{t=1}^T & \text{for classification.} \end{cases}$$

- Reduces overfitting (i.e., variance)
- Can work with any type of classifier (here focus on trees)
- Easy to parallelize (can train multiple trees in parallel)
- But loses on interpretability to single decision tree (true for all ensemble methods..)

Ensemble methods

- Bagging
- Random forests
- Boosting: Basics
- Adaboost
- Gradient boosting

Random forests

Random forests: When growing a tree on a bootstrapped (i.e. subsampled) dataset, before each split select $k \leq d$ of the d input variables at random as candidates for splitting.

When $k = d \rightarrow$ same as Bagging

When $k < d \rightarrow$ Random forests

k is a hyperparameter, tuned via cross-validation

Ensemble methods

- Bagging
- Random forests
- **Boosting: Basics**
- Adaboost
- Gradient boosting

Boosting: Idea

The boosted predictor is of the form $f_{boost}(\mathbf{x}) = \text{sign}(h(\mathbf{x}))$, where,

$$h(\mathbf{x}) = \sum_{t=1}^T \beta_t f_t(\mathbf{x}) \text{ for } \beta_t \geq 0 \text{ and } f_t \in \mathcal{F}.$$

The goal is to minimize $\ell(h(\mathbf{x}), y)$ for some loss function ℓ .

Q: We know how to find the best predictor in \mathcal{F} on some data, but how do we find the best weighted combination $h(\mathbf{x})$?

A: Minimize the loss by a *greedy approach*, i.e. find $\beta_t, f_t(\mathbf{x})$ one by one for $t = 1, \dots, T$.

Specifically, let $h_t(\mathbf{x}) = \sum_{\tau=1}^t \beta_\tau f_\tau(\mathbf{x})$. Suppose we have found $h_{t-1}(\mathbf{x})$, how do we find $\beta_t, f_t(\mathbf{x})$?

Find the $\beta_t, f_t(\mathbf{x})$ which minimizes the loss $\ell(h_t(\mathbf{x}), y)$.

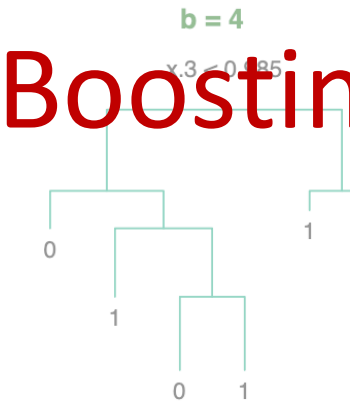
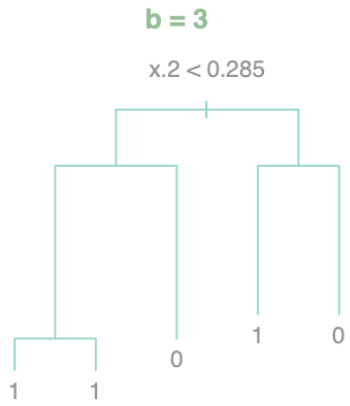
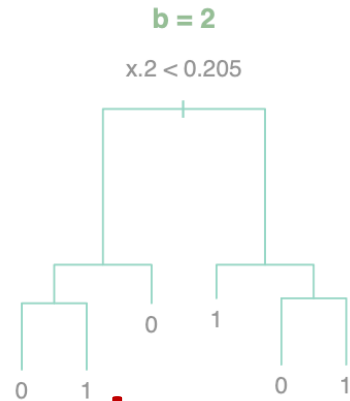
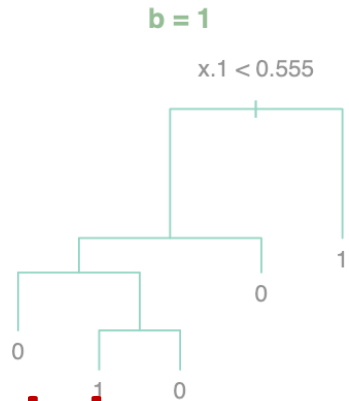
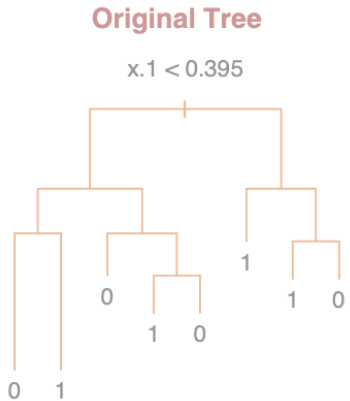
Different loss function ℓ give different boosting algorithms.

$$\ell(h(\mathbf{x}), y) = \begin{cases} (h(\mathbf{x}) - y)^2 & \rightarrow \text{Least squares boosting,} \\ \exp(-h(\mathbf{x})y) & \rightarrow \text{AdaBoost.} \end{cases}$$

Boosting: example

Email spam detection:

- given a training set like:
 - (“Want to make money fast? ...”, **spam**)
 - (“Viterbi Research Gist ...”, **not spam**)
- first obtain a classifier by applying a **base algorithm**, which can be a rather simple/weak one, like decision stumps:
 - e.g. contains the word “money” \Rightarrow spam
- **reweigh** the examples so that “**difficult**” ones get more attention
 - e.g. spam that doesn’t contain the word “money”
- obtain **another classifier** by applying the same base algorithm:
 - e.g. empty “to address” \Rightarrow spam
- repeat ...
- final classifier is the (**weighted**) **majority vote** of all weak classifiers



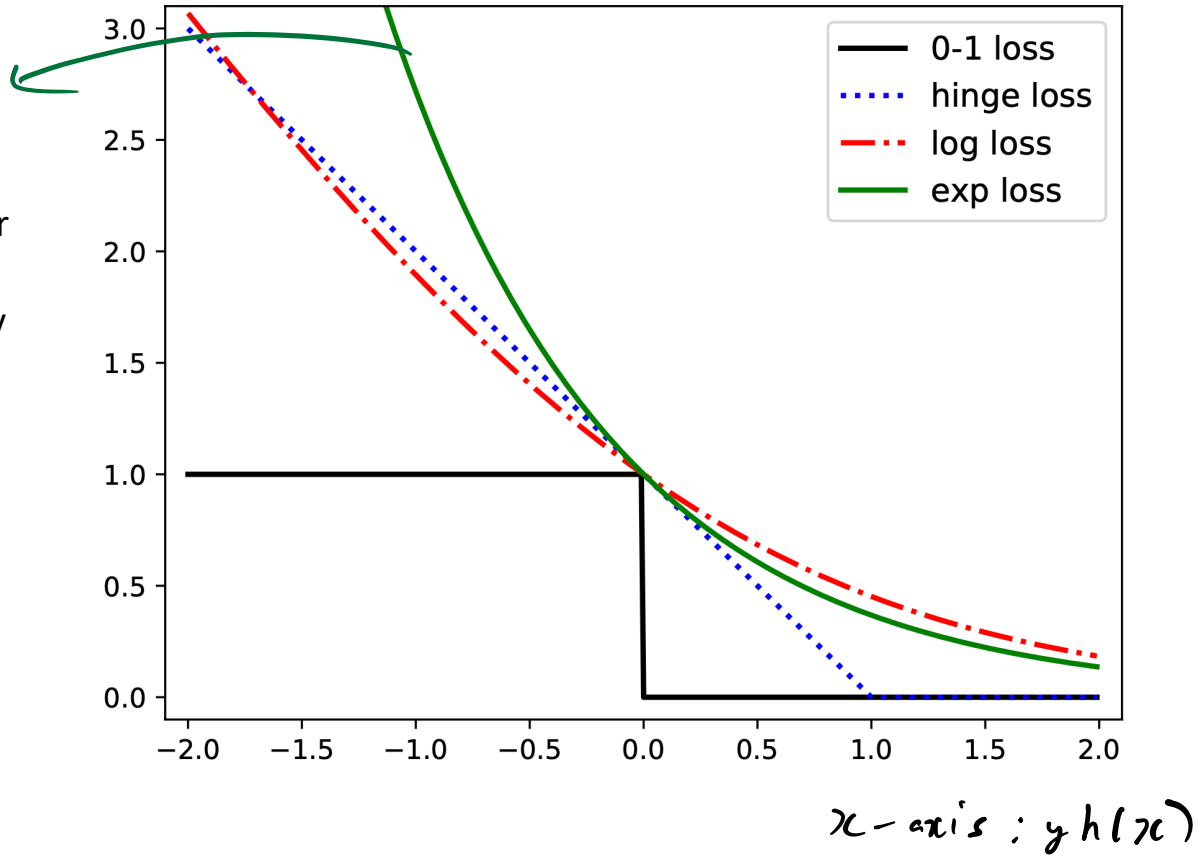
Ensemble methods: Boosting

Ensemble methods

- Bagging
- Random forests
- Boosting: Basics
- **Adaboost**
- Gradient boosting

AdaBoost

Exponential loss penalizes being far away from the decision boundary a lot more



AdaBoost

AdaBoost minimizes exponential loss by a greedy approach, that is, find $\beta_t, f_t(\mathbf{x})$ one by one for $t = 1, \dots, T$.

Recall $h_t(\mathbf{x}) = \sum_{\tau=1}^t \beta_\tau f_\tau(\mathbf{x})$. Suppose we have found h_{t-1} , *what should f_t be?* Greedily, we want to find $\beta_t, f_t(\mathbf{x})$ to minimize

$$\begin{aligned} \sum_{i=1}^n \exp(-y_i h_t(\mathbf{x}_i)) &= \sum_{i=1}^n \exp(-y_i h_{t-1}(\mathbf{x}_i)) \exp(-y_i \beta_t f_t(\mathbf{x}_i)) \\ &= \text{const}_t \cdot \sum_{i=1}^n D_t(i) \exp(-y_i \beta_t f_t(\mathbf{x}_i)) \end{aligned}$$

using $h_t = h_{t-1} + \beta_t f_t$

where the last step is by defining the weights

$$D_t(i) = \frac{\exp(-y_i h_{t-1}(\mathbf{x}_i))}{\text{const}_t}$$

const_t is a normalizing constant to make $\sum_{i=1}^n D_t(i) = 1$.

AdaBoost

So the goal becomes finding $\beta_t, f_t(\mathbf{x}) \in \mathcal{F}$ that minimize

$$\begin{aligned} & \sum_{i=1}^n D_t(i) \exp(-y_i \beta_t f_t(\mathbf{x}_i)) \quad \xrightarrow{\text{red arrow}} \quad \{f_t(\mathbf{x}_i) \in \{-1, 1\}\} \\ &= \sum_{i: y_i \neq f_t(\mathbf{x}_i)} D_t(i) e^{\beta_t} + \sum_{i: y_i = f_t(\mathbf{x}_i)} D_t(i) e^{-\beta_t} \\ &= \epsilon_t e^{\beta_t} + (1 - \epsilon_t) e^{-\beta_t} \quad (\text{where } \epsilon_t = \sum_{i: y_i \neq f_t(\mathbf{x}_i)} D_t(i) \text{ is } \textit{weighted error} \text{ of } f_t) \\ &= \epsilon_t (e^{\beta_t} - e^{-\beta_t}) + e^{-\beta_t} \end{aligned}$$

Therefore, we should find $f_t(\mathbf{x})$ to minimize its the weighted classification error ϵ_t (what we expect the base algorithm to do intuitively).

AdaBoost

When $f_t(\mathbf{x})$ (and thus ϵ_t) is fixed, we then find β_t to minimize

$$\epsilon_t(e^{\beta_t} - e^{-\beta_t}) + e^{-\beta_t}$$

Exercise: verify that the solution is given by:

$$\beta_t = \frac{1}{2} \ln \left(\frac{1 - \epsilon_t}{\epsilon_t} \right)$$

Hint: e^x is a convex function of x .

AdaBoost

How do we update the weights for the next step? The definition of $D_{t+1}(i)$ is actually recursive,

$$\begin{aligned} D_{t+1}(i) &= \frac{\exp(-y_i h_t(\mathbf{x}_i))}{\text{const}_{t+1}} \\ &= \frac{\exp(-y_i h_{t-1}(\mathbf{x}_i))}{\text{const}_{t+1}} \cdot \exp(-y_i \beta_t f_t(\mathbf{x}_i)) \\ &= \left(D_t(i) \frac{\text{const}_t}{\text{const}_{t+1}} \right) \cdot \exp(-y_i \beta_t f_t(\mathbf{x}_i)) \end{aligned}$$

$$\implies D_{t+1}(i) \propto D_t(i) \exp(-\beta_t y_i f_t(\mathbf{x}_i)) = \begin{cases} D_t(i) e^{-\beta_t} & \text{if } f_t(\mathbf{x}_i) = y_i \\ D_t(i) e^{\beta_t} & \text{else} \end{cases}$$

exponentially decrease/increase weights
based on correctly/incorrectly classified points

AdaBoost: Full algorithm

Given a training set S and a base algorithm \mathcal{A} , initialize D_1 to be uniform

For $t = 1, \dots, T$

- obtain a weak classifier $f_t(\mathbf{x}) \leftarrow \mathcal{A}(S, D_t)$
- calculate the weight β_t of $f_t(\mathbf{x})$ as

$$\beta_t = \frac{1}{2} \ln \left(\frac{1 - \epsilon_t}{\epsilon_t} \right)$$

$$(\beta_t > 0 \Leftrightarrow \epsilon_t < 0.5)$$

where $\epsilon_t = \sum_{i: f_t(\mathbf{x}_i) \neq y_i} D_t(i)$ is the weighted error of $f_t(\mathbf{x})$.

*this says that the
base algorithm does*

- update distributions

something non-trivial

$$D_{t+1}(i) \propto D_t(i) e^{-\beta_t y_i f_t(\mathbf{x}_i)} = \begin{cases} D_t(i) e^{-\beta_t} & \text{if } f_t(\mathbf{x}_i) = y_i \\ D_t(i) e^{\beta_t} & \text{else} \end{cases}$$

Output the **final classifier** $f_{boost} = \text{sgn} \left(\sum_{t=1}^T \beta_t f_t(\mathbf{x}) \right)$

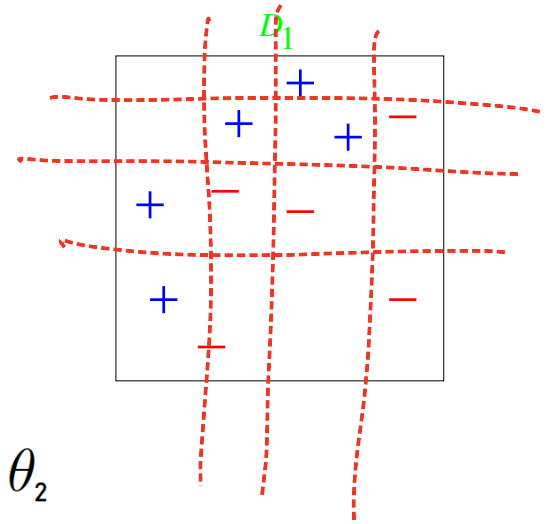
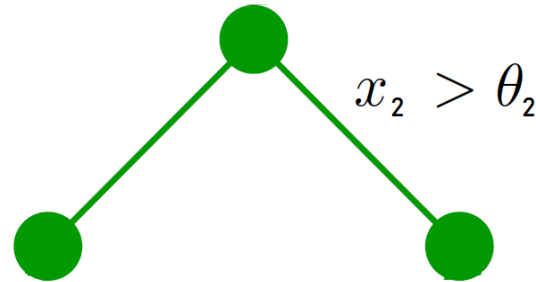
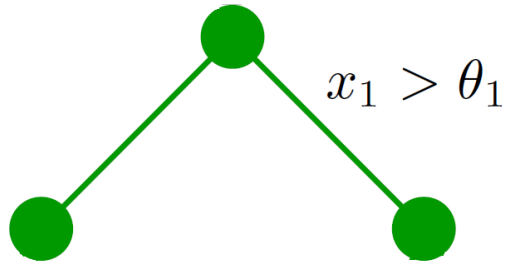
Adaboost: Example

Put more weight on difficult to classify instances and less on those already handled well
New weak learners are added sequentially that focus their training on the more difficult patterns

10 data points in \mathbb{R}^2

The size of + or - indicates the weight, which starts from uniform D_1

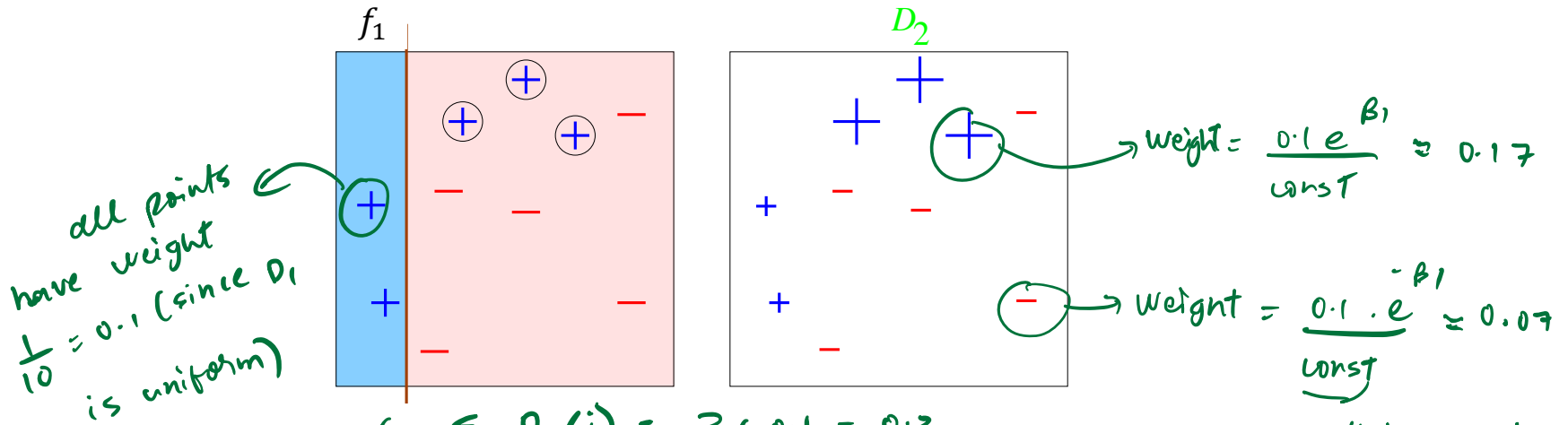
Base algorithm is decision stump:



Observe that no stump can predict very accurately for this dataset.

Adaboost: Example

Put more weight on difficult to classify instances and less on those already handled well
 New weak learners are added sequentially that focus their training on the more difficult patterns



$$\epsilon_1 = \sum_{i: f_1(x_i) \neq y_i} D_1(i) = 3 \times 0.1 = 0.3$$

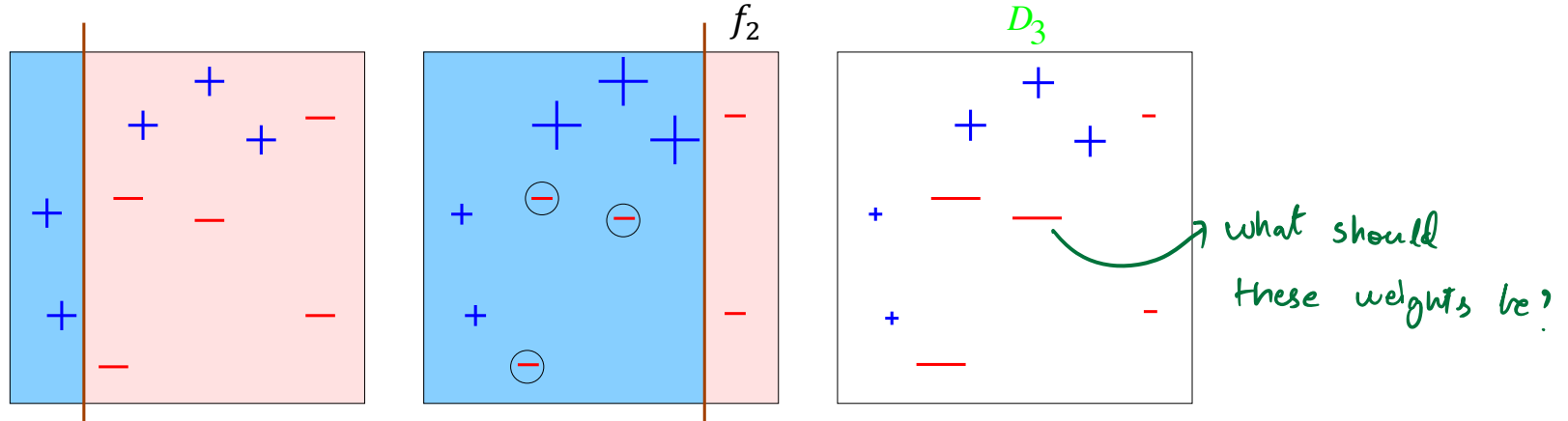
3 misclassified (circled): $\epsilon_1 = 0.3 \rightarrow \beta_1 = \frac{1}{2} \ln \left(\frac{1-\epsilon_1}{\epsilon_1} \right) \approx 0.42.$

D_2 puts more weights on these misclassified points.

normalizing const to make D_2 a distribution

Adaboost: Example

Put more weight on difficult to classify instances and less on those already handled well
New weak learners are added sequentially that focus their training on the more difficult patterns



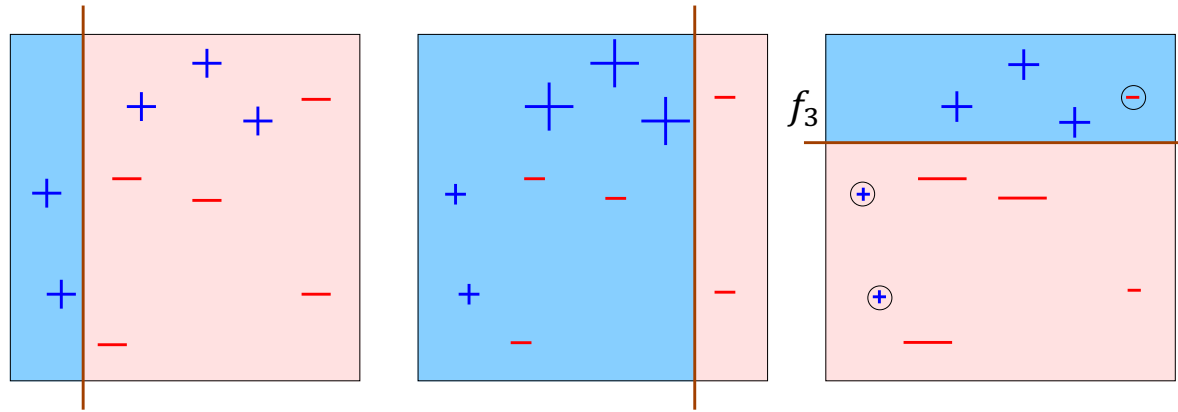
$$\epsilon_2 = \sum_{i: f_2(x_i) \neq y_i} D_2(x_i) = 0.07 + 3 = 0.21$$

3 misclassified (circled): $\epsilon_2 = 0.21 \rightarrow \beta_2 = 0.65$.

D_3 puts more weights on these misclassified points.

Adaboost: Example

Put more weight on difficult to classify instances and less on those already handled well
New weak learners are added sequentially that focus their training on the more difficult patterns

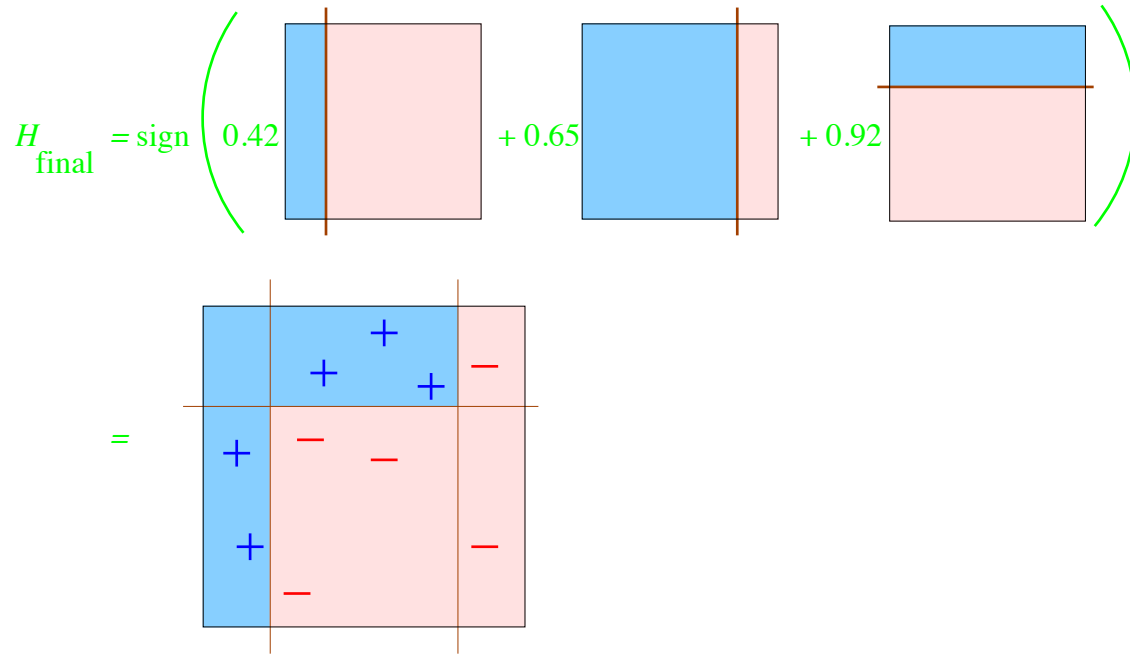


where is this coming from?

again 3 misclassified (circled): $\epsilon_3 = 0.14 \rightarrow \beta_3 = 0.92$.

Adaboost: Example

Put more weight on difficult to classify instances and less on those already handled well
New weak learners are added sequentially that focus their training on the more difficult patterns



All data points are now classified correctly, even though each weak classifier makes 3 mistakes.

Ensemble methods

- Bagging
- Random forests
- Boosting: Basics
- Adaboost
- Gradient boosting

Gradient Boosting

Recall $h_t(\mathbf{x}) = \sum_{\tau=1}^t \beta_\tau f_\tau(\mathbf{x})$. For Adaboost (exponential loss), given $h_{t-1}(\mathbf{x})$, we found what $f_t(\mathbf{x})$ should be.

Gradient boosting provides an iterative approach for general (any) loss function $\ell(h(\mathbf{x}), y)$:

- For all training datapoints (\mathbf{x}_i, y_i) find the gradient

$$r_i = - \left[\frac{\delta \ell(h(\mathbf{x}_i), y_i)}{\delta h(\mathbf{x}_i)} \right]_{h(\mathbf{x}_i) = h_{t-1}(\mathbf{x}_i)}$$

how should predictions change "locally" to reduce loss?

- Use the weak learner to find f_t which fits (\mathbf{x}_i, r_i) as well as possible:

$$f_t = \operatorname{argmin}_{f \in \mathcal{F}} \sum_{i=1}^n (r_i - f(\mathbf{x}_i))^2.$$

what should be added to reduce loss

model which tries to predict r_i or f_i

- Update $h_t(\mathbf{x}) = h_{t-1}(\mathbf{x}) + \eta f_t(\mathbf{x})$, for some step size η .

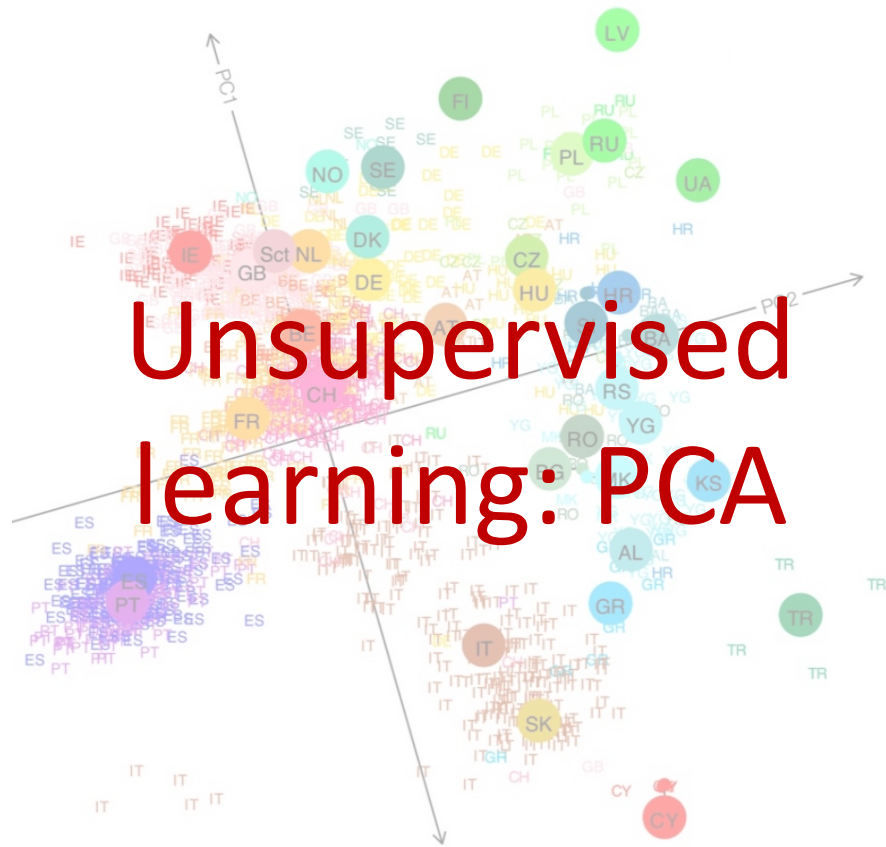
Gradient Boosting

Usually we add some regularization term to prevent overfitting (penalize the size of the tree etc.)

Gradient boosting is extremely successful!!

A variant **XGBoost** is one of the most popular algorithms for **structured data** (tables etc. with numbers and categories where each feature typically has some meaning, unlike images or text).

(for e.g. during Kaggle competitions back in 2015, 17 out of 29 winning solutions used XGBoost)



A simplistic taxonomy of ML

Supervised learning:

Aim to predict outputs of future datapoints

Unsupervised learning:

Aim to discover hidden patterns and explore data

Reinforcement learning:

Aim to make sequential decisions

Principal Component Analysis (PCA)

- Introduction
- Formalizing the problem
- How to use PCA, and examples
- Solving the PCA optimization problem
- Conclusion

Acknowledgement & further reading

Our presentation is closely based on Gregory Valiant's notes for CS168 at Stanford.

<https://web.stanford.edu/class/cs168/l/l7.pdf>

<https://web.stanford.edu/class/cs168/l/l8.pdf>

You can refer to these notes for further reading.

Also review our Linear algebra Colab notebooks:

[Part 1](#)

[Part 2](#)

Dimensionality reduction & PCA

We'll start with a simple and fundamental unsupervised learning problem: **dimensionality reduction**.

Goal: reduce the dimensionality of a dataset so that

- it is **easier to visualize and discover patterns**
- it **takes less time and space** to process for any downstream application (classification, regression, etc)
- **noise is reduced**
- ...

There are many approaches, we focus on a linear method: **Principal Component Analysis (PCA)**.

PCA: Motivation

Consider the following dataset:

- 17 features, each represents the average consumption of some food
- 4 data points, each represents some country.

What can you tell?

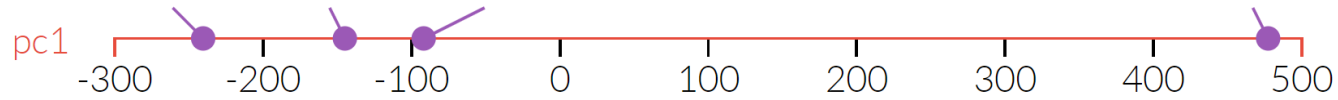
Hard to say anything looking at all these 17 features.

	England	N Ireland	Scotland	Wales
Alcoholic drinks	375	135	458	475
Beverages	57	47	53	73
Carcase meat	245	267	242	227
Cereals	1472	1494	1462	1582
Cheese	105	66	103	103
Confectionery	54	41	62	64
Fats and oils	193	209	184	235
Fish	147	93	122	160
Fresh fruit	1102	674	957	1137
Fresh potatoes	720	1033	566	874
Fresh Veg	253	143	171	265
Other meat	685	586	750	803
Other Veg	488	355	418	570
Processed potatoes	198	187	220	203
Processed Veg	360	334	337	365
Soft drinks	1374	1506	1572	1256
Sugars	156	139	147	175

Picture from [here](#)
See [this](#) for more details

PCA: Motivation

PCA can help us! The **projection of the data onto its first principal component**:

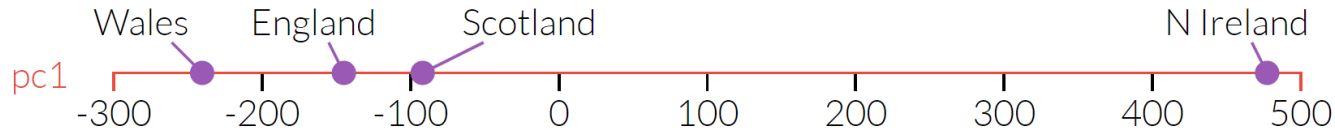


i.e. we reduce the dimensionality from 17 to just 1.

Now one data point is clearly different from the rest!

PCA: Motivation

PCA can help us! The **projection of the data onto its first principal component (PC1)**:



i.e. we reduce the dimensionality from 17 to just 1.

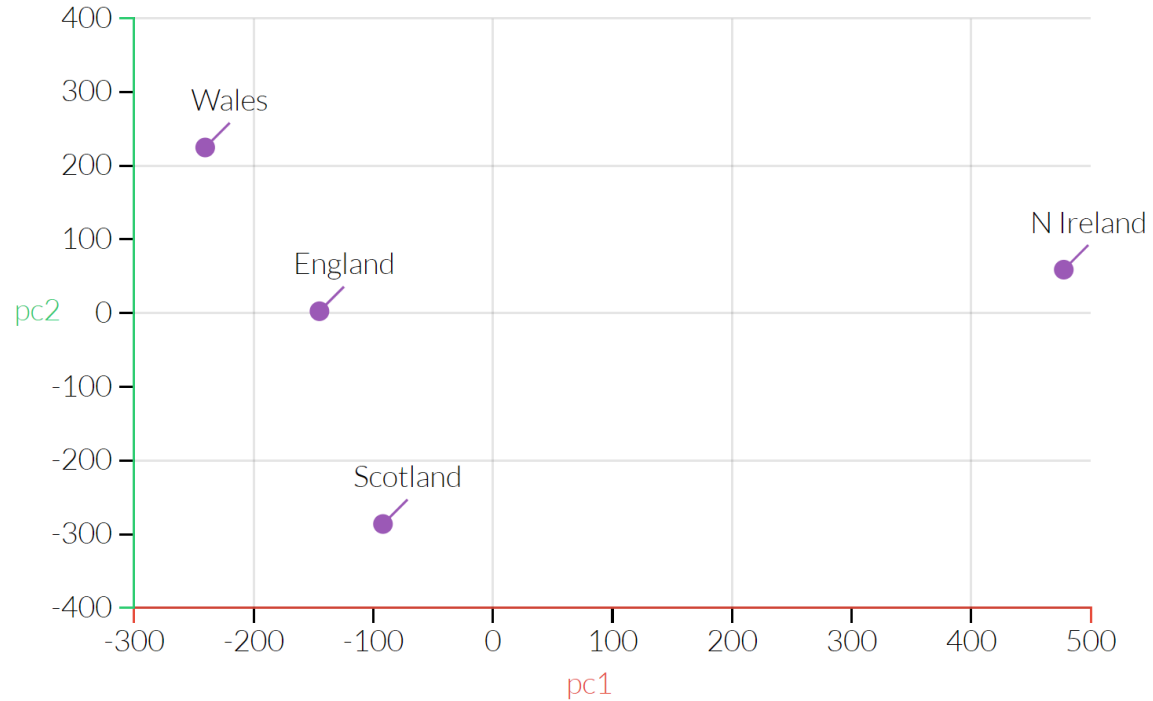
Now one data point is clearly different from the rest!

That turns out to be data from Northern Ireland,
the only country not on the island of Great Britain out of the 4 samples.

Can also interpret components: PC1 tells us that the Northern Irish eat more grams of fresh potatoes and fewer of fresh fruits and alcoholic drinks.

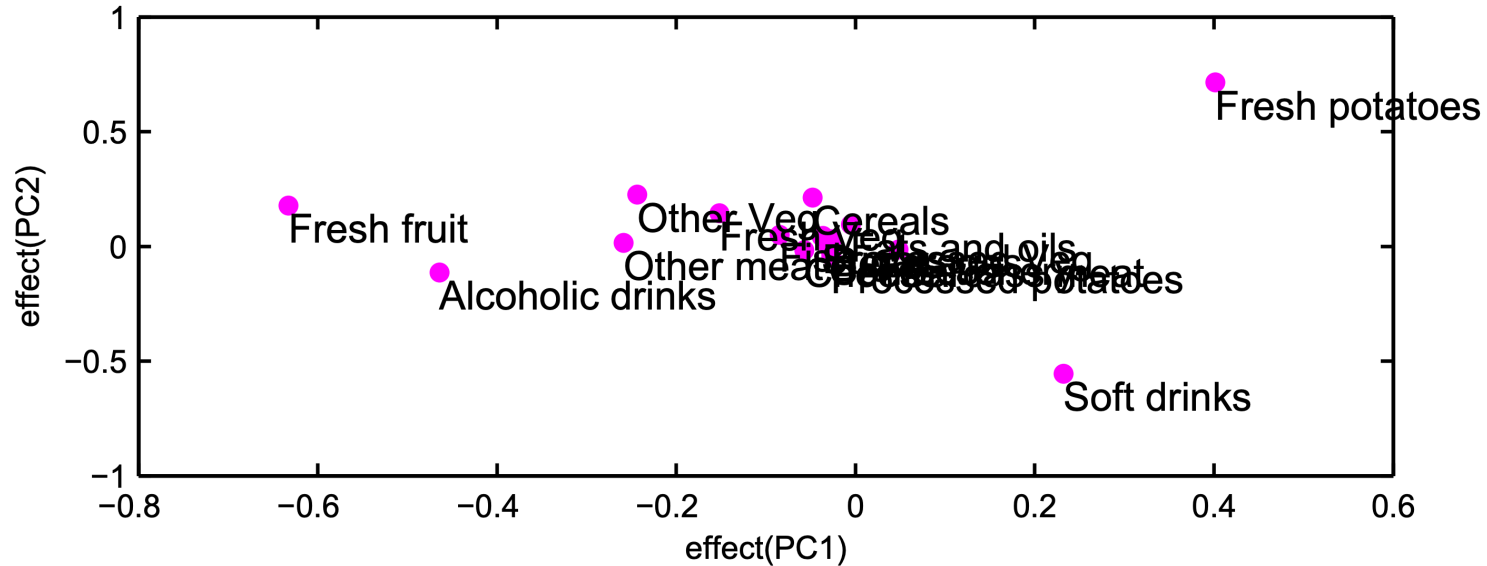
PCA: Motivation

We can find the **second (and more) principal components** of the data too:



PCA: Motivation

And the components themselves are interpretable too:



Principal Component Analysis (PCA)

- Introduction
- **Formalizing the problem**
- How to use PCA, and examples
- Solving the PCA optimization problem
- Conclusion

High-level goal

previous e.g. ·
 $n = 4$
 $d = 17$

Suppose we have a dataset of n datapoints $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n \in \mathbb{R}^d$.

The high level goal of PCA is to find a set of k principal components (PCs) / principal vectors $\mathbf{v}_1, \dots, \mathbf{v}_k \in \mathbb{R}^d$ such that for each \mathbf{x}_i ,

$$\mathbf{x}_i \approx \sum_{j=1}^k \alpha_{ij} \mathbf{v}_j$$

“principal food consumption vectors”.

food consumption vectors

for some coefficients $\alpha_{ij} \in \mathbb{R}$.

Explain the data as different linear combinations of some PCs

Preprocessing the data

- Before we apply PCA, we usually preprocess the data to center it:

$$\text{Let } \bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$$

$$\text{Centered datapoints } \tilde{\mathbf{x}}_i = \mathbf{x}_i - \bar{\mathbf{x}}$$

For the lecture, we will assume data is centered, $\sum \mathbf{x}_i = 0$

- In many applications, it is also important to scale each coordinate properly. This is especially true if the coordinates are in different units or scales.

For all $j \in [d]$, divide the j th coordinate of each point by $\sqrt{\sum_{i=1}^n \mathbf{x}_i^{j^2}}$

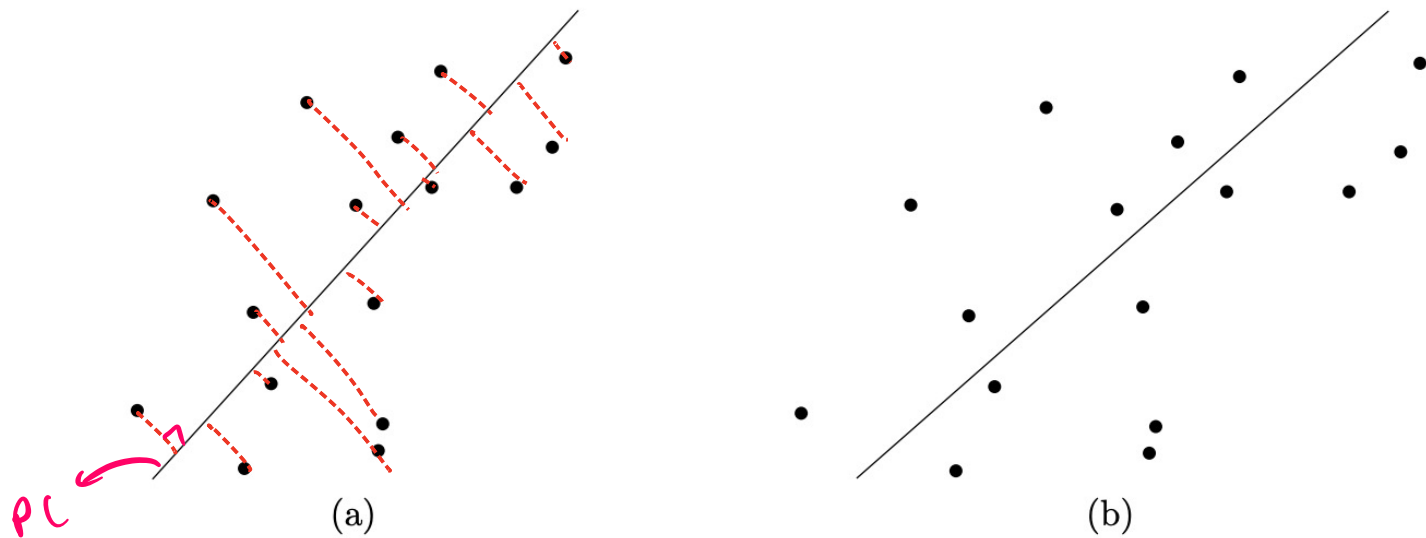


Figure 3: Scaling the x -axis yields a different best-fit line.

Objective function for PCA

Key difference from supervised learning problems:

No labels given, which means no ground-truth to measure the quality of the answer!

However, we can still write an optimization problem based on our high-level goal.

For clarity, we first discuss the special case of $k = 1$.

Optimization problem for finding the 1st principal component v_1 :

$$v_1 = \operatorname{argmin}_{v: \|v\|_2=1} \sum_{i=1}^n \left(\text{distance between } x_i \text{ \& line spanned by } v \right)^2$$

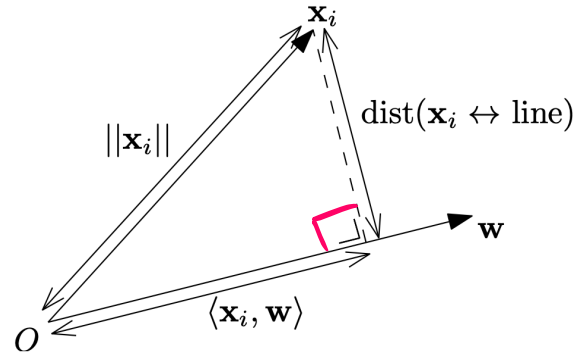


Figure 4: The geometry of the inner product with a unit length vector, w .

$$\left(\text{dist}(x_i \leftrightarrow \text{line spanned by } v) \right)^2 + \langle x_i, v \rangle^2 = \|x_i\|_2^2$$

\downarrow
 $x_i^T v$

$\|x_i\|_2^2$ is a constant, independent of choice of v

∴ original objective is equivalent to

$$\psi_1 = \arg \max_v \sum_{i=1}^n \langle x_i, v \rangle^2$$

$\psi: \|v\|_2 = 1$

An example:

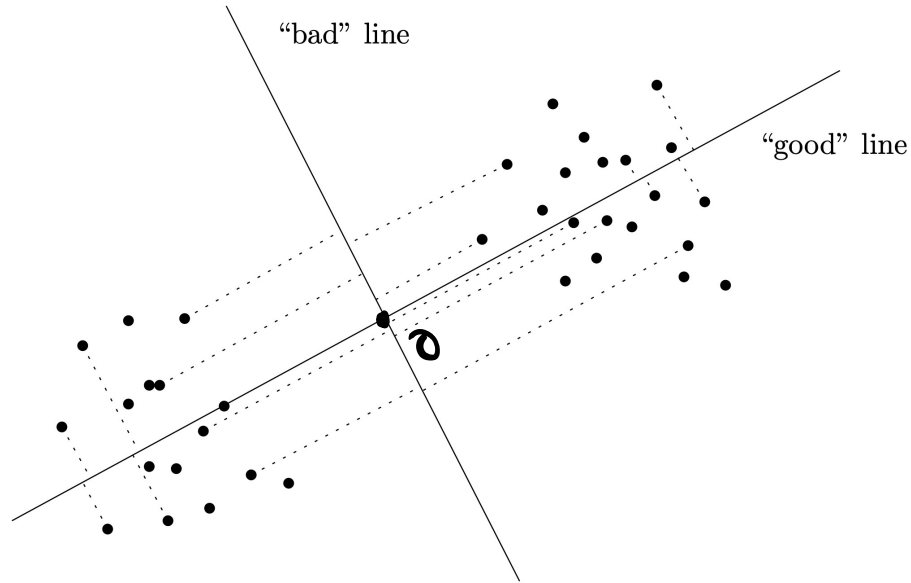



Figure 5: For the good line, the projection of the points onto the line keeps the two clusters separated, while the projection onto the bad line merges the two clusters.

Objective function for larger values of k

The generalization of the original formulation for general k is to find a k -dimensional subspace S such that the points are as close to it as possible:

$$S = \underset{k\text{-dim subspaces } S}{\operatorname{argmin}} \sum_{i=1}^n (\text{distance between } \mathbf{x}_i \text{ and subspace } S)^2$$

By the same reasoning as for $k = 1$, this is equivalent to,

$$S = \underset{k\text{-dim subspaces } S}{\operatorname{argmax}} \sum_{i=1}^n (\text{length of } \mathbf{x}_i \text{'s projection on } S)^2$$
 

It is useful to think of the subspace S as the *span* of k *orthonormal vectors* $\mathbf{v}_1, \dots, \mathbf{v}_k \in \mathbb{R}^d$.

Recall, vectors $\mathbf{v}_1, \dots, \mathbf{v}_k$ are orthonormal if:

1. $\|\mathbf{v}_i\|_2 = 1 \quad \forall i \in [k]$
2. $\langle \mathbf{v}_i, \mathbf{v}_j \rangle = 0 \quad \forall i \neq j$

e.g.
standard basis vectors
 $\mathbf{e}_1 = (1, 0, 0, \dots, 0)$
 $\mathbf{e}_2 = (0, 1, 0, \dots, 0)$
 \vdots

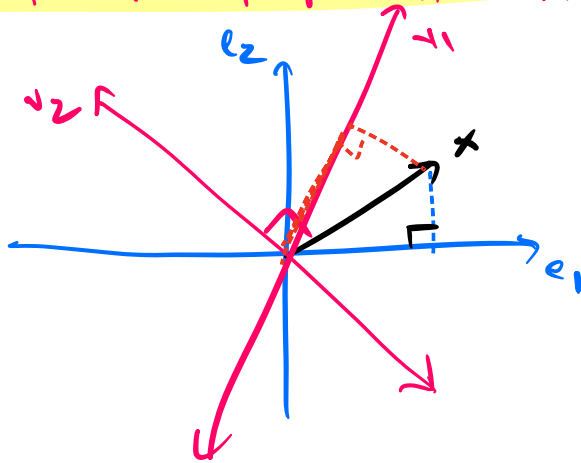
Definition: Span of a collection $v_1, \dots, v_k \in \mathbb{R}^d$ is all their linear combinations $\left\{ \sum_{j=1}^k \lambda_j v_j : \lambda_1, \dots, \lambda_k \in \mathbb{R} \right\}$

Example,

- $k = 1$, span is line through the origin.
- $k = 2$, if v_1, v_2 are linearly independent, the span is a plane through the origin, and so on.

Fact about orthogonal vectors

$$\left(\text{length of } x \text{'s projection on } \text{span}(v_1, \dots, v_k) \right)^2 = \sum_{j=1}^k \langle x, v_j \rangle^2 \quad (2)$$



Combining ① & ②

Formal problem solved by PCA:

Given $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$ and a parameter $k \geq 1$, compute orthonormal vectors $\mathbf{v}_1, \dots, \mathbf{v}_k \in \mathbb{R}^d$ to maximize,

$$\sum_{i=1}^n \sum_{j=1}^k \langle \mathbf{x}_i, \mathbf{v}_j \rangle^2.$$

Equivalent view:

- Pick \mathbf{v}_1 to be the variance maximizing direction.
- Pick \mathbf{v}_2 to be the variance maximizing direction, orthogonal to \mathbf{v}_1 .
- Pick \mathbf{v}_3 to be the variance maximizing direction, orthogonal to \mathbf{v}_1 and \mathbf{v}_2 , and so on.

Principal Component Analysis (PCA)

- Introduction
- Formalizing the problem
- **How to use PCA, and examples**
- Solving the PCA optimization problem
- Conclusion

Using PCA for data compression and visualization

Input: n datapoints $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n \in \mathbb{R}^d$, #components k we want

Step 1 Perform PCA to get top k principal components $\mathbf{v}_1, \dots, \mathbf{v}_k \in \mathbb{R}^d$.

Step 2 For each datapoint \mathbf{x}_i , define its “ \mathbf{v}_1 -coordinate” as $\langle \mathbf{x}_i, \mathbf{v}_1 \rangle$, its “ \mathbf{v}_2 -coordinate” as $\langle \mathbf{x}_i, \mathbf{v}_2 \rangle$. Therefore we associate k coordinates to each datapoint \mathbf{x}_i , where the j -th coordinate denotes the extent to which \mathbf{x}_i points in the direction of \mathbf{v}_j .

Step 3 We now have a new “compressed” dataset where each datapoint is k -dimensional. For visualization, we can plot the point \mathbf{x}_i in \mathbb{R}^k as the point $(\langle \mathbf{x}_i, \mathbf{v}_1 \rangle, \langle \mathbf{x}_i, \mathbf{v}_2 \rangle, \dots, \langle \mathbf{x}_i, \mathbf{v}_k \rangle)$.

going back to high-level goal:

$$\mathbf{x}_i \approx \sum_{j=1}^k \langle \mathbf{x}_i, \mathbf{v}_j \rangle \mathbf{v}_j .$$

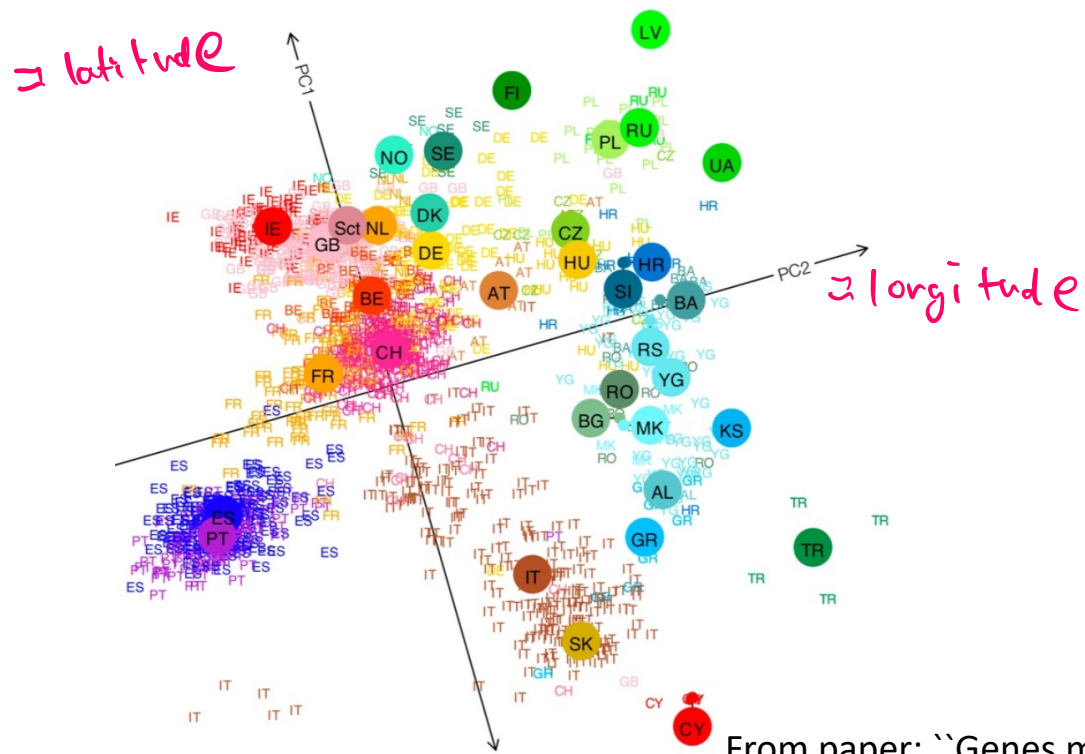
Visualization example: Do Genomes Encode Geography?

Dataset: genomes of 1,387 Europeans (each individual's genotype at 200,000 locations in the genome)

$n = 1387, d \approx 200,000$

Project the datapoints onto top 2 PCs

Plot shown below; looks remarkably like the map of Europe!



From paper: "Genes mirror geography within Europe" Novembre et al., Nature'08

Compression example: Eigenfaces

Dataset: 256x256 ($\approx 65\text{K}$ pixels) dimensional images of about 2500 faces, all framed similarly
 $n = 2500, d \approx 65,000$

We can represent each image with high accuracy using only 100-150 principal components!

The principal components (called *eigenfaces* here) are themselves interpretable too!

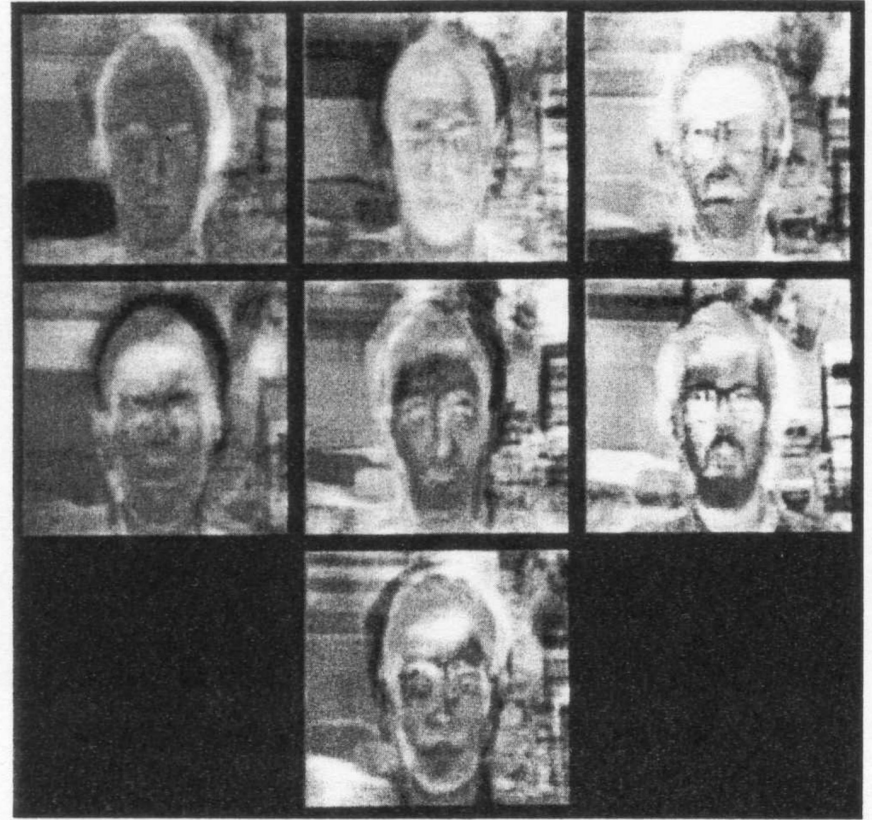


Figure 2. Seven of the eigenfaces calculated from the input images of Figure 1.

Principal Component Analysis (PCA)

- Introduction
- Formalizing the problem
- How to use PCA, and examples
- Solving the PCA optimization problem
- Conclusion

How to solve the PCA optimization problem?

Consider $k=1$,

$$v_1 = \operatorname{arg\,max}_{v: \|v\|_2=1} \sum_{i=1}^n \langle x_i, v \rangle^2$$

$$X = \begin{bmatrix} \text{---} x_1^T \text{---} \\ \text{---} x_2^T \text{---} \\ \vdots \\ \text{---} x_n^T \text{---} \end{bmatrix} \in \mathbb{R}^{n \times d}$$

\therefore for any $v \in \mathbb{R}^d$,

$$Xv = \begin{bmatrix} x_1^T v \\ x_2^T v \\ \vdots \\ x_n^T v \end{bmatrix} \in \mathbb{R}^n$$

$$\begin{aligned}
 \therefore \sum_{i=1}^n (x_i^T v)^2 &= \|Xv\|_2^2 \\
 &= (Xv)^T (Xv) \\
 &= v^T \underbrace{X^T X}_A v
 \end{aligned}$$

\therefore for $A = X^T X \in \mathbb{R}^{d+d}$

$$v_1 = \operatorname{argmax}_v v^T A v$$

$$v : \|v\|_2 = 1$$

$X^T X$: covariance of data (data is centered)

$$A = \begin{bmatrix} | & | & & | \\ x_1 & x_2 & \dots & x_n \\ | & | & & | \end{bmatrix} \begin{bmatrix} \text{--- } x_1^T \text{ ---} \\ \text{--- } x_2^T \text{ ---} \\ \vdots \\ \text{--- } x_n^T \text{ ---} \end{bmatrix}$$

A_{11} : variance of 1st coordinate

$$(A_{11} = \sum_{i=1}^n x_{i,1}^2)$$

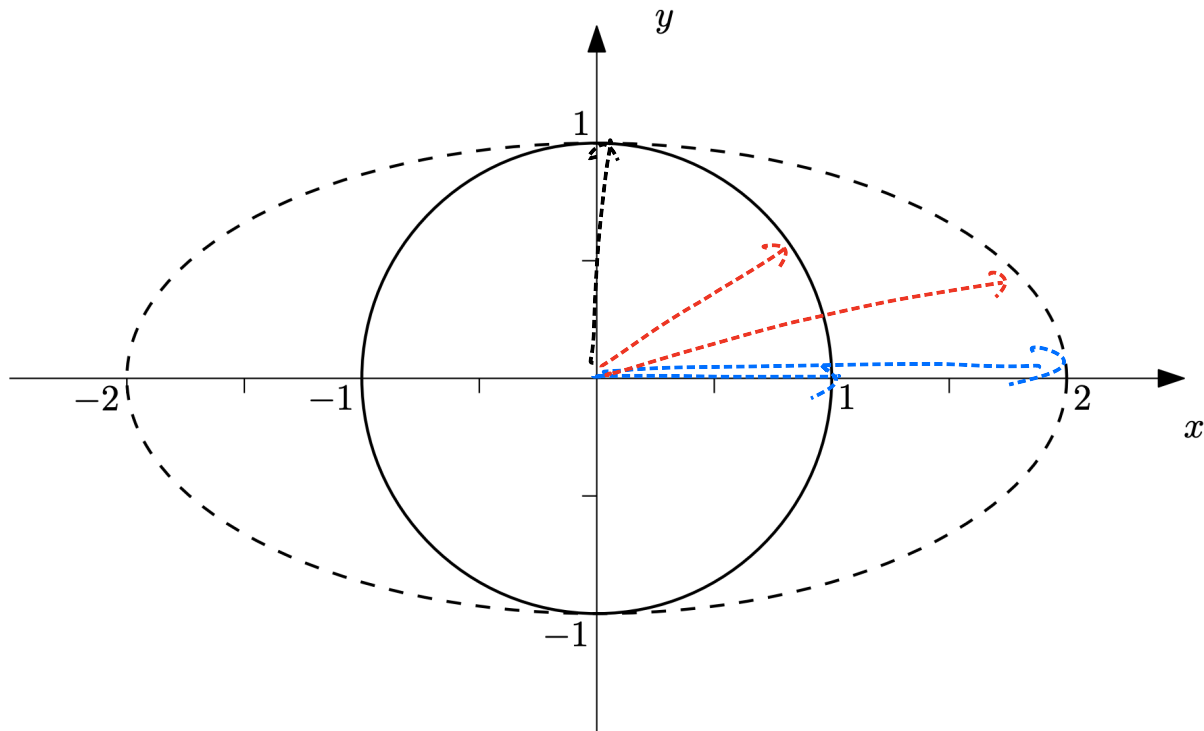
The diagonal case

Let's solve $\operatorname{argmax}_{\mathbf{v}: \|\mathbf{v}\|_2=1} \mathbf{v}^\top \mathbf{A} \mathbf{v}$ for the special case where \mathbf{A} is a diagonal matrix.

$$\mathbf{A} = \begin{pmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \lambda_d \end{pmatrix} \quad \text{where } \lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_d \geq 0$$

Any $d \times d$ matrix \mathbf{A} can be thought of as a function that maps points in \mathbb{R}^d back to points in \mathbb{R}^d : $\mathbf{v} \mapsto \mathbf{A} \mathbf{v}$.

The matrix $\begin{pmatrix} 2 & 0 \\ 0 & 1 \end{pmatrix}$ maps (x, y) to $(2x, y)$:



Points on circle $\{(x, y) : x^2 + y^2 = 1\}$ are mapped to the ellipse $\{(x, y) : (\frac{x}{2})^2 + y^2 = 1\}$.

So what direction v should maximize $v^T A v$ for diagonal A ?

It should be the direction of maximum stretch: (assuming $\lambda_1 > \lambda_2 \geq \dots \geq \lambda_d \geq 0$)
 $v = e_1$ (where $e_i = (1, 0, \dots, 0)$ is 1st standard basis vector)

Proof:

$$v^T A v = v^T (A v) = (v_1 \dots v_d) \begin{pmatrix} \lambda_1 v_1 \\ \lambda_2 v_2 \\ \vdots \\ \lambda_d v_d \end{pmatrix} = \sum_{i=1}^d v_i^2 \lambda_i$$

Since v is a unit vector, $\sum_{i=1}^d v_i^2 = 1$

\therefore Since λ_1 is largest, to maximize set $v_1 = 1$
 $v_i = 0 \quad \forall i > 1.$

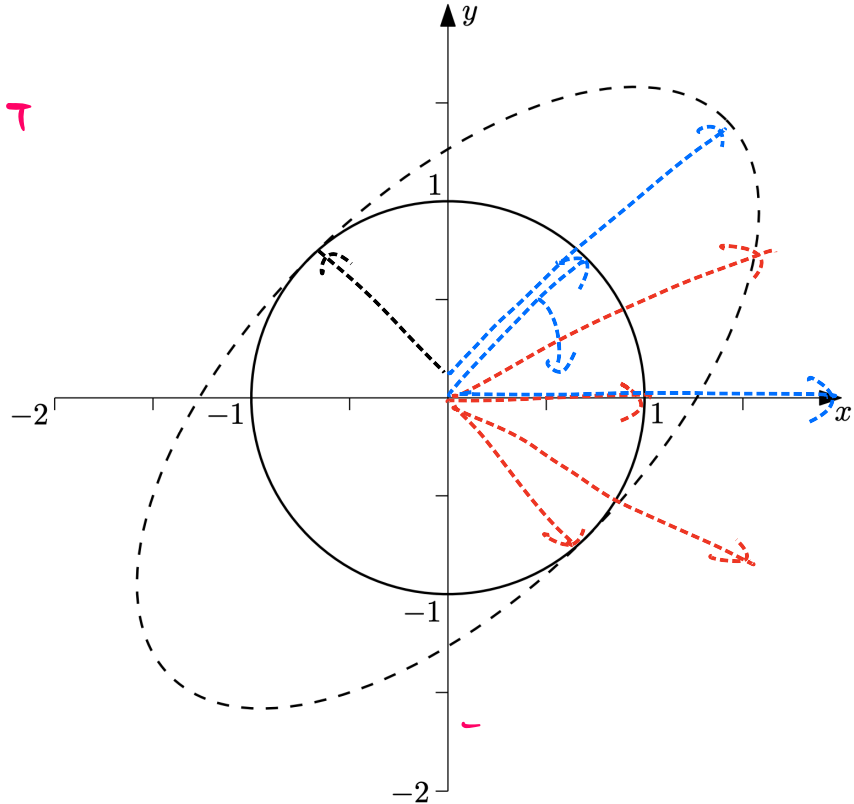
$\therefore \vec{v} = e_1$

Diagonals in disguise

Consider

$$\begin{aligned} \mathbf{A} &= \begin{pmatrix} \frac{3}{2} & \frac{1}{2} \\ \frac{1}{2} & \frac{3}{2} \end{pmatrix} \xrightarrow{\mathbf{Q}} \\ &= \begin{pmatrix} \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{pmatrix} \cdot \begin{pmatrix} 2 & 0 \\ 0 & 1 \end{pmatrix} \xrightarrow{\mathbf{D}} \begin{pmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ -\frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{pmatrix} \xrightarrow{\mathbf{Q}^T} \end{aligned}$$

\mathbf{A} still does nothing other than stretch out different orthogonal axes, possibly with these axes being a “rotated version” of the original ones.



The previous figure, rotated 45 degrees.

How do we formalize the concept of a rotation in high dimensions as a matrix operation?

Answer: Orthogonal matrix (also called orthonormal matrix).

An orthonormal matrix is a matrix Q s.t. for all columns

$q_1, \dots, q_d,$

$$\|q_i\|_2 = 1 \quad \forall i$$

$$q_i^T q_j = 0 \quad \forall i \neq j$$

Key properties:

$$\textcircled{1} \quad Q^T Q = I \quad (Q^{-1} = Q^T)$$

$$\textcircled{2} \quad \|Qv\|_2^2 = \|v\|_2^2$$

$\textcircled{3}$ If Q is orthogonal, Q^T is also orthogonal

$$\textcircled{1} \quad \begin{pmatrix} -q_1^T- \\ -q_2^T- \\ \vdots \\ -q_d^T- \end{pmatrix} \begin{pmatrix} | & | & | \\ q_1 & q_2 & \dots & q_d \\ | & | & | \end{pmatrix} \\ = \begin{pmatrix} 1 & & 0 \\ & \ddots & \\ 0 & & 1 \end{pmatrix}$$

$$\textcircled{2} \quad \|Qv\|_2^2 = (Qv)^T Qv = v^T \underbrace{Q^T Q}_I v \\ = v^T v = \|v\|_2^2$$

$\textcircled{3}$ Exercise!

Recall that we want to find $\mathbf{v}_1 = \operatorname{argmax}_{\mathbf{v}: \|\mathbf{v}\|_2=1} \mathbf{v}^T \mathbf{A} \mathbf{v}$.

Now consider \mathbf{A} that can be written as $\mathbf{A} = \mathbf{Q} \mathbf{D} \mathbf{Q}^T$ for an orthogonal matrix \mathbf{Q} and diagonal matrix \mathbf{D} with diagonal entries $\lambda_1 \geq \lambda_2 \geq \lambda_3 \geq \dots \lambda_d \geq 0$.

$$\mathbf{A} = \begin{array}{|c|} \hline \begin{array}{c} | \quad | \quad | \\ \mathbf{Q}_1 \ \mathbf{Q}_2 \ \dots \ \mathbf{Q}_d \\ | \quad | \quad | \end{array} \\ \hline \end{array} \quad \begin{array}{|c|} \hline \begin{array}{c} \lambda_1 \quad \quad \quad 0 \\ \quad \lambda_2 \quad \quad \quad \vdots \\ \quad \quad \quad \vdots \quad \quad \quad 0 \\ \quad \quad \quad 0 \quad \quad \quad \lambda_d \end{array} \\ \hline \end{array} \quad \begin{array}{|c|} \hline \begin{array}{c} \text{--- } \mathbf{Q}_1^T \text{ ---} \\ \text{--- } \mathbf{Q}_2^T \text{ ---} \\ \vdots \\ \text{--- } \mathbf{Q}_d^T \text{ ---} \end{array} \\ \hline \end{array} \quad \begin{array}{c} \left(\begin{array}{c} \mathbf{Q} \mathbf{e}_1 \\ \mathbf{v} \end{array} \right) \end{array}$$

only rotate, so doesn't affect length

What is the direction which gets stretched the maximum?

(Informal answer) The maximum possible stretch by \mathbf{D} is λ_1 . The direction of maximum stretch under \mathbf{D} is \mathbf{e}_1 . Therefore, direction of maximum stretch under $\mathbf{D} \mathbf{Q}^T$ is \mathbf{v} s.t. $\mathbf{Q}^T \mathbf{v} = \mathbf{e}_1 \implies \mathbf{v} = (\mathbf{Q}^T)^{-1} \mathbf{e}_1 = \mathbf{Q} \mathbf{e}_1$.

Claim: For $A = QDQ^T$

$$\text{argmax}_v v^T A v = Qe_1$$

$$v: \|v\|_2 = 1$$

Proof: for $v_i = Qe_1$

$$\begin{aligned} v_i^T A v_i &= (Qe_1)^T A (Qe_1) = e_1^T \underbrace{Q^T}_{I} Q D \underbrace{Q^T}_{I} Q e_1 \\ &= e_1^T D e_1 \\ &= d_{11} \end{aligned}$$

$$\text{Now for } v^T A v = \left(\begin{array}{c} \text{---} \\ v^T \end{array} \right) \left(\begin{array}{c} \square \\ Q \end{array} \right) \left(\begin{array}{c} \square \\ \underbrace{D}_{d_{11}} \end{array} \right) \left(\begin{array}{c} \square \\ Q^T \\ \square \\ e \end{array} \right)$$

Q & Q^T preserve length $\therefore Q^T Q$ (& $v^T Q$) are unit vectors.

$$\therefore v^T A v \leq d_{11}$$

$\therefore v_i = Qe_1$ maximizes $v^T A v$. \bullet

General covariance matrices

Consider any covariance $A = X^T X$

Linear algebra fact: Any symmetric A can be written as $A = Q D Q^T$ for orthogonal matrix Q & diagonal matrix D .

If $A = X^T X$, A is symmetric & D always has non-negative entries. Why?

$$v^T A v = v^T X^T X v = \|X v\|_2^2 \geq 0$$

If $D_{ii} < 0$, then $v = Q e_i$, $v^T Q D Q^T v = D_{ii} < 0$

When $k = 1$, the solution to $\operatorname{argmax}_{v: \|v\|_2=1} v^T A v$ is the first column of Q , where $A = X^T X = Q D Q^T$ with Q orthogonal and D diagonal with sorted entries.

General values of k

What is the solution to the PCA objective for general values of k ?

$$\sum_{i=1}^n \sum_{j=1}^k \langle \mathbf{x}_i, \mathbf{v}_j \rangle^2$$

Solution: Pick the first k columns of \mathbf{Q} , where the covariance $\mathbf{X}^T \mathbf{X} = \mathbf{Q} \mathbf{D} \mathbf{Q}^T$ with \mathbf{Q} orthogonal and \mathbf{D} diagonal with sorted entries.

Since \mathbf{Q} is orthogonal, the first k columns of \mathbf{Q} are orthogonal vectors. These are called the top k principal components (PCs).

Eigenvalues & eigenvectors

How to compute the top k columns of Q in the decomposition $X^T X = Q D Q^T$?

Solution: Eigenvalue decomposition!

Definition: An eigenvector of matrix A is a vector v that is stretched in the same direction by A , i.e.,

$$A v = \lambda v$$

for some $\lambda \in \mathbb{R}$. Here λ is the eigenvalue corresponding to the eigenvector v .

Eigenvectors: axes of stretch in geometric intuition

Eigenvalues: stretch factors

When we write $A = X^T X$ as $A = Q D Q^T$

→ rows of Q^T (columns of Q) are eigenvectors of A

→ diagonal entries of D are corresponding eigenvalues.

Proof:

i th column of Q is given by $Q e_i$ $\left(\begin{array}{|c|c|c|c|} \hline | & | & | & | \\ \hline Q_1 & Q_2 & \dots & Q_d \\ \hline | & | & | & | \\ \hline \end{array} \right)$

$$A(Q e_i) = Q D \underbrace{Q^T Q}_{I} e_i = Q \underbrace{D e_i}_{\lambda_i e_i} = Q \lambda_i e_i = \lambda_i (Q e_i)$$

∴ i th column of Q is an eigenvector with eigenvalue λ_i .

PCA boils down to computing the k eigenvectors of the covariance matrix $X^T X$ that have the largest eigenvalues.

Principal Component Analysis (PCA)

- Introduction
- Formalizing the problem
- How to use PCA, and examples
- Solving the PCA optimization problem
- **Conclusion**

How many PCs to use?

For visualization, we usually choose k to be small and just pick the first few principal components.

In other applications such as compression, it is a good idea to plot the eigenvalues and see. A lot of data is close to being low rank, so the eigenvalues may decay and become small.

We can also choose the threshold based on how much variance we want to capture. Suppose we want to capture 90% of the variance in the data. Then we can pick k such that i.e.

$$\frac{\sum_{j=1}^k \lambda_j}{\sum_{j=1}^d \lambda_j} \geq 90\%$$

where $\lambda_1 \geq \dots \geq \lambda_d$ are sorted eigenvalues.

Note: $\sum_{j=1}^d \lambda_j = \text{trace}(\mathbf{X}^T \mathbf{X})$, so no need to actually find all eigenvalues.
sum of diagonal entries of $\mathbf{X}^T \mathbf{X}$

When and why does PCA fail?

1. Data is not properly scaled/normalized.
2. Non-orthogonal structure in data: PCs are forced to be orthogonal, and there may not be too many orthogonal components in the data which are all interpretable.
3. Non-linear structure in data.

