

CSCI-567: Machine Learning

Prof. Victor Adamchik

U of Southern California

July 22, 2020

Your model is only as good as your data.

July 22, 2020 1 / 48

Outline

1 Principal Component Analysis (PCA)

- PCA
- Kernel PCA

2 Boosting

3 Problem Solving

Outline

1 Principal Component Analysis (PCA)

2 Boosting

3 Problem Solving

July 22, 2020 2 / 48

Dimensionality reduction

Dimensionality reduction is yet another important unsupervised learning problem.

Goal: reduce the dimensionality of a dataset so

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

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

July 22, 2020 3 / 48

July 22, 2020 4 / 48

Example

Consider the following dataset:

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

	Alcoholic drinks	Beverages	Carcase meat	Cereals	Cheese	Confectionery	Fats and oils	Fish	Fresh fruit	Fresh potatoes	Fresh Veg	Other meat	Other Veg	Processed potatoes	Processed Veg	Soft drinks	Sugars
Alcoholic drinks	375	47	135	458	53	475	73										
Beverages	57	245	47	267	41	242	227										
Carcase meat	245	105	267	1494	66	1462	1582										
Cereals	1472	54	267	1494	41	1462	1582										
Cheese	105	54	66	103	103	62	64										
Confectionery	54	193	209	184	184	235											
Fats and oils	193	147	93	122	122	160											
Fish	147	102	674	957	957	1137											
Fresh fruit	102	720	1033	566	566	874											
Fresh potatoes	720	253	143	171	171	265											
Fresh Veg	253	685	586	750	750	803											
Other meat	685	488	355	418	418	570											
Other Veg	488	198	187	220	220	203											
Processed potatoes	198	360	334	337	337	365											
Processed Veg	360	1374	1509	1572	1572	1256											
Soft drinks	1374	156	139	147	147	175											
Sugars	156																

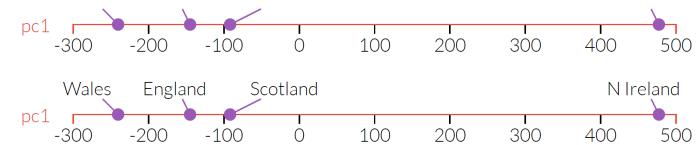
What can you tell?

Hard to say anything looking at all these 17 features.

July 22, 2020 5 / 48

Example

PCA can help us! The first principal component of this dataset:



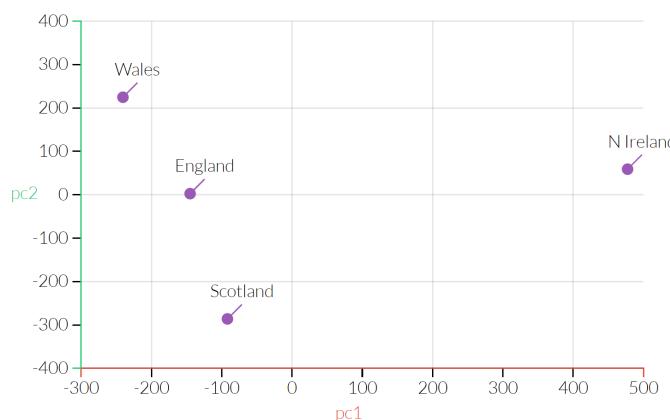
i.e. we reduce the dimensionality from 17 to just 1.

That turns out to be data from Northern Ireland,

July 22, 2020 6 / 48

Example

PCA can find the second principal component of the data too:



High level idea

How does PCA find these principal components (PC)?



This is in fact the direction with the most variance, i.e. the direction where the data is most spread out.

July 22, 2020 7 / 48

July 22, 2020 8 / 48

Finding the first PC

More formally, we want to find a direction $\mathbf{v} \in \mathbb{R}^D$ with $\|\mathbf{v}\|_2 = 1$, so that the **projection of the dataset on this direction has the most variance**, i.e.

$$\max_{\mathbf{v}: \|\mathbf{v}\|_2=1} \sum_{n=1}^N \left(\mathbf{x}_n^\top \mathbf{v} - \frac{1}{N} \sum_m \mathbf{x}_m^\top \mathbf{v} \right)^2$$

- $\mathbf{x}_n^\top \mathbf{v}$ is exactly the **projection of \mathbf{x}_n onto the direction \mathbf{v}**
- if we **pre-center the data**, i.e. let $\mathbf{x}_n \leftarrow \mathbf{x}_n - \frac{1}{N} \sum_m \mathbf{x}_m$, then the objective simply becomes

$$\max_{\mathbf{v}} \sum_{n=1}^N (\mathbf{x}_n^\top \mathbf{v})^2 = \max_{\mathbf{v}} \sum_{n=1}^N (\mathbf{v}^\top \mathbf{x}_n \mathbf{x}_n^\top \mathbf{v}) = \max_{\mathbf{v}} \mathbf{v}^\top \left(\sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top \right) \mathbf{v}$$

July 22, 2020 9 / 48

Finding the other PCs

If v_1 is the first PC, then the second PC is found via

$$\max_{\mathbf{v}_2: \|\mathbf{v}_2\|_2=1, \mathbf{v}_1^\top \mathbf{v}_2=0} \mathbf{v}_2^\top (\mathbf{X}^\top \mathbf{X}) \mathbf{v}_2$$

i.e. the direction that maximizes the variance among all other dimensions.

This is just the **second top eigenvector of the covariance matrix!**

Conclusion: the d -th principal component is the d -th eigenvector (sorted by the eigenvalue from largest to smallest).

Finding the first PC

With $\mathbf{X} \in \mathbb{R}^{N \times D}$ being the data matrix, we want

$$\max_{\mathbf{v}: \|\mathbf{v}\|_2=1} \mathbf{v}^\top (\mathbf{X}^\top \mathbf{X}) \mathbf{v}$$

Let the max be $\lambda > 0$, then

$$\mathbf{v}^\top (\mathbf{X}^\top \mathbf{X}) \mathbf{v} = \lambda = \lambda \mathbf{v}^\top \mathbf{v} = \mathbf{v}^\top (\lambda \mathbf{v})$$

It follows,

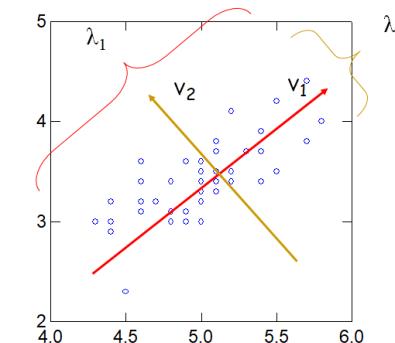
$$\mathbf{v}^\top (\mathbf{X}^\top \mathbf{X} \mathbf{v} - \lambda \mathbf{v}) = 0$$

Conclusion: the first PC is the top eigenvector of the covariance matrix.

July 22, 2020 10 / 48

PCA

PCA eigenvectors:



Input: a dataset represented as \mathbf{X} , #components p we want

Step 1 Center the data by subtracting the mean

Step 2 Find the top p (unit norm) eigenvectors of the covariance matrix $\mathbf{X}^T \mathbf{X}$, denote it by $\mathbf{V} \in \mathbb{R}^{D \times p}$ (each column is an eigenvector).

Step 3 Construct the new compressed dataset $\mathbf{X}\mathbf{V} \in \mathbb{R}^{N \times p}$

July 22, 2020 13 / 48

How many PCs do we want?

One common rule: pick p large enough so it **covers about 90% of the spectrum**, i.e.

$$\frac{\sum_{d=1}^p \lambda_d}{\sum_{d=1}^D \lambda_d} \geq 90\%$$

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

Note: $\sum_{d=1}^D \lambda_d = \text{Tr}(\mathbf{X}^T \mathbf{X})$, so no need to actually find all eigenvalues.

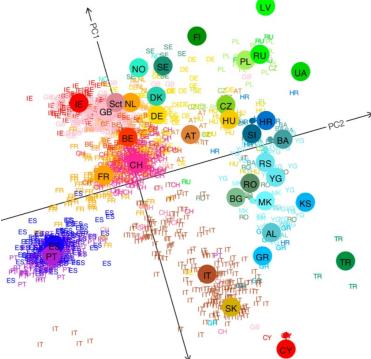
For **visualization**, also often pick $p = 1$ or $p = 2$.

July 22, 2020 14 / 48

Another visualization example

A famous study of **genetic map**

- dataset: genomes of 1,387 Europeans
- First 2 PCs shown below; *looks remarkably like the geographic map*

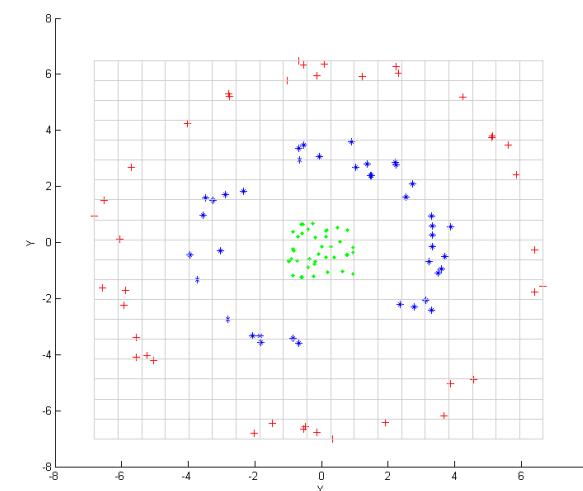


July 22, 2020 15 / 48

Does PCA always work?

picture from Wikipedia

PCA is a **linear method** (recall the new dataset is $\mathbf{X}\mathbf{V}$), it does not do much when **every direction has similar variance**.



July 22, 2020 16 / 48

KPCA: high level idea

Similar to learning a linear classifier, when we encounter such data, *we can apply kernel methods.*

Kernel PCA (KPCA):

- first map the data to a more complicated space via $\phi : \mathbb{R}^D \rightarrow \mathbb{R}^M$
- then apply regular PCA to reduce the dimensionality

Sounds a bit counter-intuitive, but the key is this gives a *nonlinear method*.

How to implement KPCA efficiently without actually working in \mathbb{R}^M ?

One issue: scaling

Should we scale α s.t $\|\alpha\|_2 = 1$?

No. Recall we want $v = \Phi^T \alpha$ to have unit L2 norm:

$$v^T v = \alpha^T \Phi \Phi^T \alpha = \alpha^T \lambda \alpha = \lambda \|\alpha\|_2^2 = 1$$

In fact we need to *scale α so that its L2 norm is $1/\sqrt{\lambda}$* , where λ it's the corresponding eigenvalue.

KPCA: finding the PCs

Suppose $v \in \mathbb{R}^M$ is the first PC for the nonlinearly-transformed data $\Phi \in \mathbb{R}^{N \times M}$ (centered). Then

$$v = \frac{1}{\lambda} \Phi^T \Phi v = \Phi^T \alpha$$

for some $\alpha \in \mathbb{R}^N$, i.e. v is a linear combination of data.

Plugging that v into $\Phi^T \Phi v = \lambda v$ gives

$$\Phi^T \Phi \Phi^T \alpha = \lambda \Phi^T \alpha$$

and thus with the Gram matrix $K = \Phi \Phi^T$,

$$\Phi^T K \alpha = \lambda \Phi^T \alpha$$

So α is an eigenvector of K !

Conclusion: KPCA is just finding top eigenvectors of the Gram matrix.

Another issue: centering

Should we still pre-center our data set X ?

No. Centering X does not mean Φ is centered!

Remember all we need is Gram matrix. What is the Gram matrix \bar{K} after Φ is centered?

Let $E \in \mathbb{R}^{N \times N}$ be the matrix with all entries being $\frac{1}{N}$, then $E\Phi$ is a matrix where each row is the mean of data

$$\begin{aligned}\bar{K} &= (\Phi - E\Phi)(\Phi - E\Phi)^T \\ &= (\Phi - E\Phi)(\Phi^T - \Phi^T E) \\ &= \Phi \Phi^T - E\Phi \Phi^T - \Phi \Phi^T E + E\Phi \Phi^T E \\ &= K - EK - KE + EKE\end{aligned}$$

Input: a dataset \mathbf{X} , #components p we want, a Kernel function k .

Step 1 Compute the Gram matrix \mathbf{K} and the centered Gram matrix

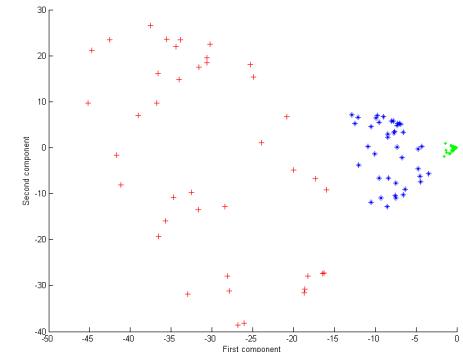
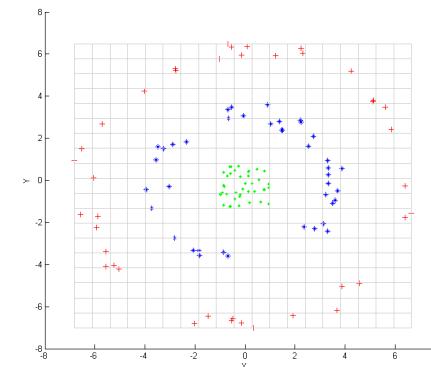
$$\bar{\mathbf{K}} = \mathbf{K} - \mathbf{E}\mathbf{K} - \mathbf{K}\mathbf{E} + \mathbf{E}\mathbf{K}\mathbf{E}$$

Step 2 Find the top p eigenvectors of $\bar{\mathbf{K}}$ with the appropriate scaling, denote it by $\mathbf{A} \in \mathbb{R}^{N \times p}$. Each column is an eigenvector.

Step 3 Construct the new dataset \mathbf{KA}

Example

Applying kernel $k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^T \mathbf{x}' + 1)^2$:



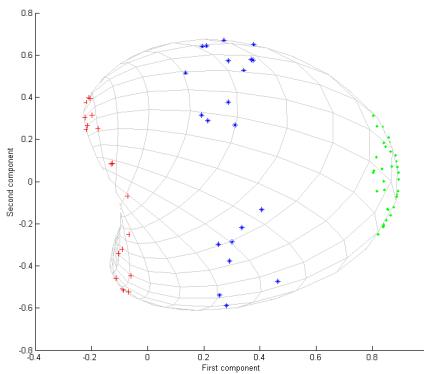
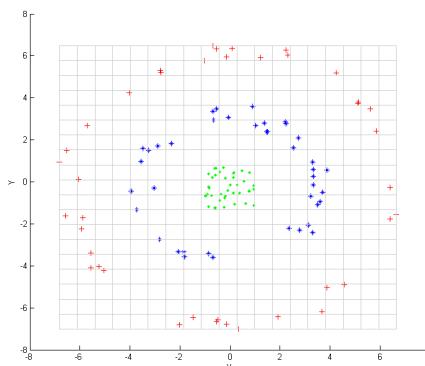
July 22, 2020 21 / 48

July 22, 2020 22 / 48

Example

picture from Wikipedia

Applying Gaussian kernel $k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x}-\mathbf{x}'\|^2}{2\sigma^2}\right)$:



July 22, 2020 23 / 48

Denoising via PCA

Original data



Data corrupted with Gaussian noise



Result after linear PCA



Result after kernel PCA, Gaussian kernel



July 22, 2020 24 / 48

1 Principal Component Analysis (PCA)

2 Boosting

- Examples
- AdaBoost
- Derivation of AdaBoost

3 Problem Solving

A simple 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
- reweight 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

Boosting

- is a **meta-algorithm**, which takes a base algorithm (classification, regression, ranking, etc) as input and **boosts** its accuracy
- main idea: combine **weak "rules of thumb"** (e.g. 51% accuracy) to form a **highly accurate predictor** (e.g. 99% accuracy)
- works very well in practice (especially in combination with trees)
- often is **resistant to overfitting**
- has strong theoretical guarantees

We again focus on **binary classification**.

The base algorithm

A **base algorithm** \mathcal{A} (also called weak learning algorithm/oracle) takes a **training set S** **weighted by D** as input, and outputs classifier $h \leftarrow \mathcal{A}(S, D)$

- this can be **any off-the-shelf classification algorithm** (e.g. decision trees, logistic regression, neural nets, etc)
- many algorithms can deal with a **weighted training set** (e.g. for algorithm that minimizes some loss, we can simply **replace "total loss"** by "**weighted total loss**")

Boosting Algorithms

Given:

- a training set S
- a base algorithm \mathcal{A}

Two things to specify a boosting algorithm:

- how to **reweight** the examples?
- how to **combine** all the weak classifiers?

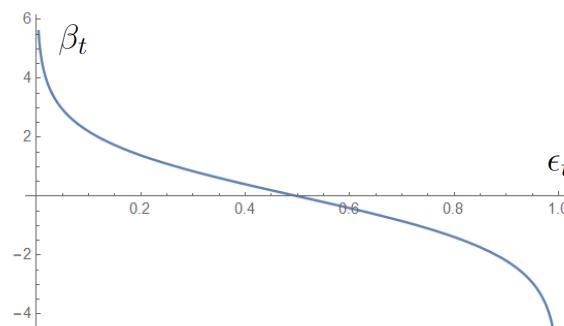
AdaBoost is one of the most successful boosting algorithms.

July 22, 2020 29 / 48

The Betas

Calculate the **importance** of h_t as

$$\beta_t = \frac{1}{2} \ln \left(\frac{1 - \epsilon_t}{\epsilon_t} \right) \quad (\beta_t > 0 \Leftrightarrow \epsilon_t < 0.5)$$



July 22, 2020 31 / 48

The AdaBoost Algorithm, 1990

Given N samples $\{\mathbf{x}_n, y_n\}$, where $y_n \in \{+1, -1\}$, and a base algorithm \mathcal{A} .

Initialize $D_1 = \frac{1}{N}$ to be **uniform**.

For $t = 1, \dots, T$

- Train a weak classifier $h_t \leftarrow \mathcal{A}(S, D_t)$ based on the current weight $D_t(n)$, by minimizing the **weighted** classification error

$$\epsilon_t = \sum_n D_t(n) \mathbb{I}[y_n \neq h_t(\mathbf{x}_n)]$$

- Calculate the **importance** of h_t as

$$\beta_t = \frac{1}{2} \ln \left(\frac{1 - \epsilon_t}{\epsilon_t} \right) \quad (\beta_t > 0 \Leftrightarrow \epsilon_t < 0.5)$$

July 22, 2020 30 / 48

The AdaBoost Algorithm

For $t = 1, \dots, T$

- Train a weak classifier $h_t \leftarrow \mathcal{A}(S, D_t)$
- Calculate β_t
- **Update weights**

$$D_{t+1}(n) = D_t(n) e^{-\beta_t y_n h_t(\mathbf{x}_n)} = \begin{cases} D_t(n) e^{-\beta_t} & \text{if } h_t(\mathbf{x}_n) = y_n \\ D_t(n) e^{\beta_t} & \text{else} \end{cases}$$

and normalize them such that $D_{t+1}(n) = \frac{D_{t+1}(n)}{\sum_n D_{t+1}(n)}$.

Output the **final classifier**:

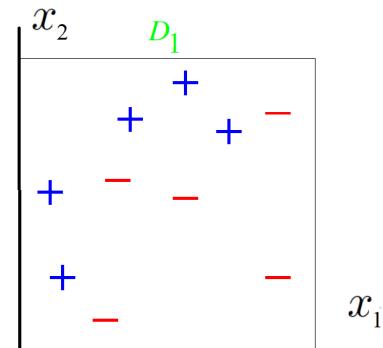
$$H(\mathbf{x}) = \text{sign} \left(\sum_{t=1}^T \beta_t h_t(\mathbf{x}) \right)$$

July 22, 2020 32 / 48

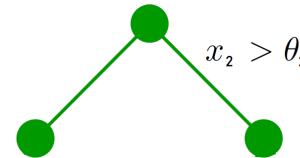
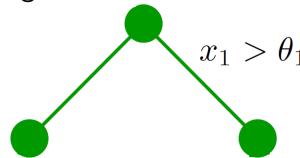
Example

10 data points in \mathbb{R}^2

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

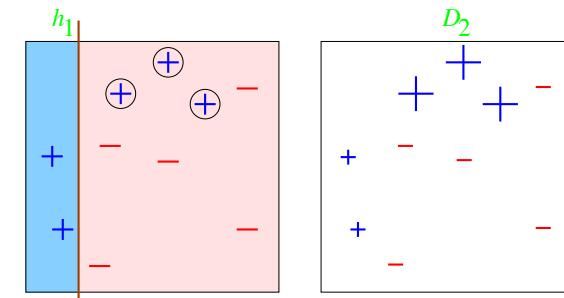


Base algorithm is decision stump:



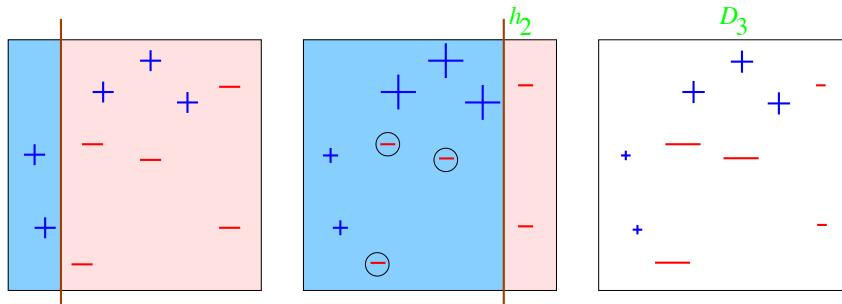
Observe that *no stump can predict very accurately for this dataset*

Round 1: $t = 1$



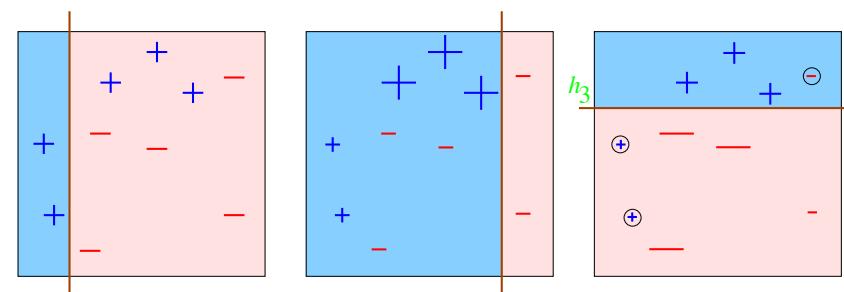
- 3 misclassified (circled): $\epsilon_1 = 0.3 \rightarrow \beta_1 = \frac{1}{2} \ln \left(\frac{1-\epsilon_t}{\epsilon_t} \right) \approx 0.42$.
- D_2 puts more weights on those examples

Round 2: $t = 2$



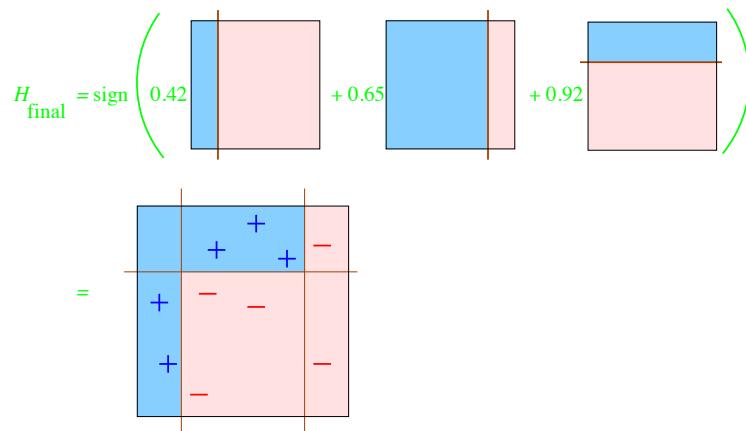
- 3 misclassified (circled): $\epsilon_2 = 0.21 \rightarrow \beta_2 = 0.65$.
- D_3 puts more weights on those examples

Round 3: $t = 3$



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

Final classifier: combining 3 classifiers



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

Why AdaBoost works?

In fact, AdaBoost also follows the general framework of minimizing some surrogate loss.

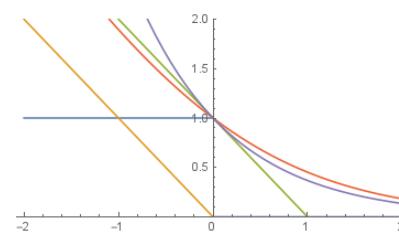
Step 1: the model that AdaBoost considers is

$$\left\{ \text{sgn}(f(\cdot)) \mid f(\cdot) = \sum_{t=1}^T \beta_t h_t(\cdot) \text{ for some } \beta_t \geq 0 \text{ and } h_t \in \mathcal{H} \right\}$$

where \mathcal{H} is the set of models considered by the base algorithm

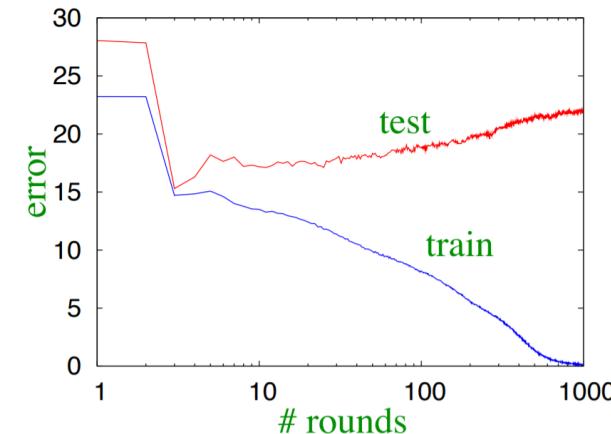
Step 2: the loss that AdaBoost minimizes is the exponential loss

$$\sum_{n=1}^N \exp(-y_n f(\mathbf{x}_n))$$



Overfitting

When T is large, the model is very complicated and overfitting can happen



(boosting “stumps” on heart-disease dataset)

However, very often AdaBoost is resistant to overfitting

Greedy minimization

Step 3: the way that AdaBoost minimizes exponential loss is by a greedy approach, that is, find β_t, h_t one by one for $t = 1, \dots, T$.

Specifically, let $f_t = \sum_{\tau=1}^t \beta_\tau h_\tau$. Suppose we have found f_{t-1} , what should f_t be?

$$f_t = \sum_{\tau=1}^{t-1} \beta_\tau h_\tau + \beta_t h_t = f_{t-1} + \beta_t h_t.$$

Greedily, we want to find β_t, h_t to minimize

$$\sum_{n=1}^N \exp(-y_n f_t(\mathbf{x}_n)) = \sum_{n=1}^N \exp(-y_n f_{t-1}(\mathbf{x}_n)) \exp(-y_n \beta_t h_t(\mathbf{x}_n))$$

Next, we use the definition of weights (slide 32).

Greedy minimization

Claim:

$$\exp(-y_n f_{t-1}(\mathbf{x}_n)) \propto D_t(n)$$

Proof.

$$\begin{aligned} D_t(n) &\propto D_{t-1}(n) \exp(-y_n \beta_{t-1} h_{t-1}(\mathbf{x}_n)) \\ &\propto D_{t-2}(n) \exp(-y_n \beta_{t-2} h_{t-2}(\mathbf{x}_n)) \exp(-y_n \beta_{t-1} h_{t-1}(\mathbf{x}_n)) \\ &\propto D_1 \exp(-y_n \beta_1 h_1(\mathbf{x}_n) - \dots - y_n \beta_{t-1} h_{t-1}(\mathbf{x}_n)) \\ &\propto \exp(-y_n f_{t-1}(\mathbf{x}_n)) \end{aligned}$$

Remark. All weights $D_t(n)$ are normalized: $\sum_n D_t(n) = 1$.

Minimizing the weighted classification error

Thus, we would want to choose $h_t(\mathbf{x}_n)$ such that

$$h_t^*(\mathbf{x}) = \operatorname{argmin}_{h_t(\mathbf{x})} \epsilon_t = \sum_{n:y_n \neq h_t(\mathbf{x}_n)} D_t(n)$$

This is exactly the first step of the AdaBoost algorithm on slide 30 — *train a weak classifier based on the current weight $D_t(n)$* .

Greedy minimization

So the goal becomes finding $\beta_t \geq 0, h_t \in \mathcal{H}$ that minimize

$$\operatorname{argmin}_{\beta_t, h_t} \sum_{n=1}^N \exp(-y_n f_t(\mathbf{x}_n)) = \operatorname{argmin}_{\beta_t, h_t} \sum_{n=1}^N D_t(n) \exp(-y_n \beta_t h_t(\mathbf{x}_n))$$

We decompose the weighted loss function into two parts

$$\begin{aligned} &\sum_{n=1}^N D_t(n) \exp(-y_n \beta_t h_t(\mathbf{x}_n)) \\ &= \sum_{n:y_n \neq h_t(\mathbf{x}_n)} D_t(n) e^{\beta_t} + \sum_{n:y_n = h_t(\mathbf{x}_n)} D_t(n) e^{-\beta_t} \\ &= \epsilon_t e^{\beta_t} + (1 - \epsilon_t) e^{-\beta_t} \quad (\text{recall } \epsilon_t = \sum_{n:y_n \neq h_t(\mathbf{x}_n)} D_t(n)) \\ &= \epsilon_t (e^{\beta_t} - e^{-\beta_t}) + e^{-\beta_t} \end{aligned}$$

We find h_t by minimizing the weighted classification error ϵ_t .

Greedy minimization

When h_t (and thus ϵ_t) is fixed, we then find β_t to minimize

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

We take derivative with respect to β_t , set it to zero, and derive the optimal β_t as

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

which is precisely the second step of the AdaBoost algorithm on slide 30.

Verify the solution β_t^* .

Updating the weights

Now that we have improved our classifier into

$$f_t(\mathbf{x}_n) = f_{t-1}(\mathbf{x}_n) + \beta_t^* h_t^*(\mathbf{x}_n)$$

At the t -th iteration, we will need to compute the weights for the above classifier, which is,

$$\begin{aligned} D_{t+1}(n) &\propto e^{-y_n f(\mathbf{x}_n)} = e^{-y_n [f_{t-1}(\mathbf{x}) + \beta_t^* h_t^*(\mathbf{x}_n)]} \\ &\propto D_t(n) e^{-y_n \beta_t^* h_t^*(\mathbf{x}_n)} = \begin{cases} D_t(n) e^{\beta_t^*} & \text{if } y_n \neq h_t^*(\mathbf{x}_n) \\ D_t(n) e^{-\beta_t^*} & \text{if } y_n = h_t^*(\mathbf{x}_n) \end{cases} \end{aligned}$$

which is precisely the last step of the AdaBoost algorithm on slide 32.

Summary for boosting

Key idea of boosting is to **combine weak predictors into a strong one**.

There are many boosting algorithms; AdaBoost is the most classic one.

AdaBoost is **greedily minimizing the exponential loss**.

AdaBoost tends to **not overfit**.

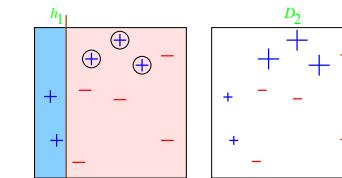
Remarks

Note that the AdaBoost algorithm itself never specifies how we would get $h_t^*(\mathbf{x})$ as long as it minimizes the weighted classification error

$$\epsilon_t = \sum_n D_t(n) \mathbb{I}[y_n \neq h_t^*(\mathbf{x}_n)]$$

In this aspect, the AdaBoost algorithm is a meta-algorithm and can be used with any classifier where we can do the above.

Ex. How do we choose the decision stump classifier given the weights at the second round of the following distribution?



We can simply enumerate all possible ways of putting vertical and horizontal lines to separate the data points into two classes and find the one with the smallest weighted classification error!

Outline

1 Principal Component Analysis (PCA)

2 Boosting

3 Problem Solving

Problem 1

Find the second PCA value:

$$\max(v_2^T X^T X v_2),$$

subject to

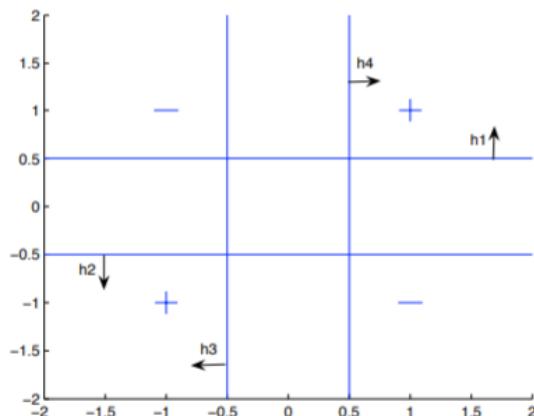
$$v_2^T v_2 = 1$$

$$v_2^T v_1 = 0$$

Solution

Problem 2

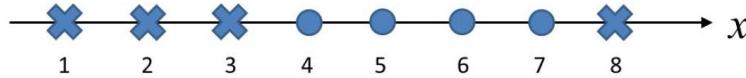
Consider the four binary classifiers below. The arrow means that the corresponding classifier classifies every data point in that direction as +. Prove that there are no weights β_1, \dots, β_4 , that make the ensemble $\sum \beta_i h_i$ classifier consistent with the data.



Solution

Problem 3

Imagine running AdaBoost with a 1-dimensional training set of 8 examples as shown



Circles mean $y = +1$ and crosses mean $y = -1$. The number under each example is its x coordinate. The base classifier set H consists of all decision stumps such that

$$h_i(x) = \begin{cases} s, & \text{if } x > b \\ -s, & \text{otherwise} \end{cases}$$

1. Run AdaBoost for two rounds and compute β_1 and β_2 .
2. Compute the training error of the final classifier H .

Solution
