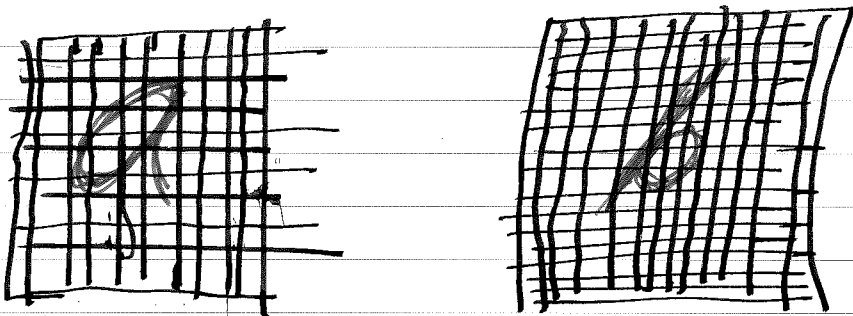


Some Pattern Recognition & Overview

PCA was an example of unsupervised learning algorithm in that it segregates data without feedback or training examples.

Consider the Character Recognition problem



Character 'a' & 'b' are images represented as an array of pixels the i^{th} pixel having value x_i , say, where $x_i = 0$ means white & $x_i = 1$ means black. Together the image is a vector $\hat{x} = [x_1 \dots x_d]^T$

Want a function $f: \mathbb{R}^d \rightarrow \{1, 2\}$ to be a character recognizer. Really f is "classifying" by mapping \mathbb{R}^d into \mathcal{C}_1 or $\mathcal{C}_2 \leftarrow$ class 1 or class 2.

We imagine we have a training set of classified images

$$\{(\hat{x}^{(i)}, y^{(i)})\}_{i=1}^m$$

An image might be a 256×256

object so $d = 65,536$. If colors are stored as

8-bit ints that's $2^{8 \cdot 65,536} \approx 10^{158,000}$ possible

images! That's huge as there are only $\sim 10^{80}$ particles

available in the universe.

On the other hand our training set might only have a few thousand images.

The classifier, f , is then desired to have the ability to generalize, namely to correctly classify unseen-before examples.

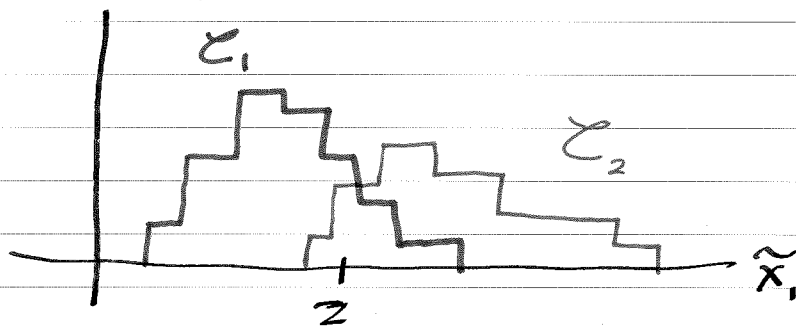
Note that the curse of dimensionality is acting against us here.

One way to try to fight dimensionality is to create new variables, called features which combine input variables and reduce the size of dimension. Design of features is important.

For example, maybe $\tilde{x}_1 = \frac{\text{height of character}}{\text{width of character}}$

\tilde{x}_1 then combines tons of input variables (pixel values) into a single feature variable which, wairly, is more informative.

Consider the Histogram



How should a new character with $\tilde{x}_1 = 2$ be classified?

Since there is significant overlap a great # of misclassifications will occur if we use a simple thresholding, namely,
let $x^* = \{ \tilde{x}_1 \mid C_1 \text{ wins} = C_2 \text{ wins} \}$

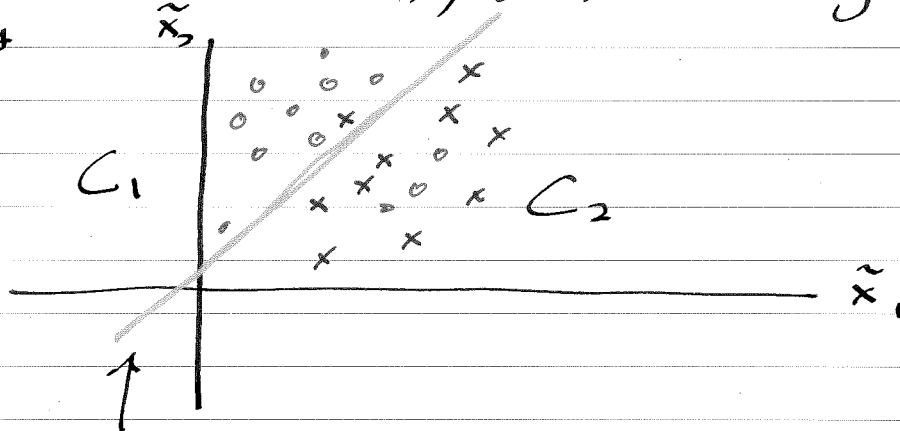
$$\text{Then } f_{\text{thres}}(\tilde{x}) = \begin{cases} C_1 & \text{if } \forall \tilde{x}_1(\tilde{x}) < x^* \\ C_2 & \text{if } \tilde{x}_1(\tilde{x}) \geq x^* \end{cases}$$

But the size of overlap indicates there will likely be many misclassifications. since the map

$\tilde{x}_1: \mathbb{R}^d \rightarrow \mathbb{R}$ represents some loss of information.

Perhaps there's a second feature \tilde{x}_2 (whose details don't matter!)

and we can then represent the training set in \mathbb{R}^2



called a "decision boundary" and separates $(\tilde{x}_1, \tilde{x}_2)$ space into the 2 class regions.

Now $f_{dec}(\tilde{x}) = \begin{cases} C_1 & \text{if } \tilde{x} \text{ is above boundary} \\ C_2 & \text{else} \end{cases}$

In general with mappings

$$\mathbb{R}^d \ni \tilde{x} \longrightarrow y_k \in \{0, 1, 2, 3, \dots, N\}$$

$N = \#$ of classifiers as if classifying all letters would have $N = 26$.

In general $f(\vec{x}) = y_k(\vec{x}; \vec{\theta})$

$\vec{\theta}$ = vector of model parameters or often called "weights".

Determining $\vec{\theta}$ based on data set is called "learning" or "training".

Then Classification problems are when $\text{Range}(f)$ is a discrete set.

When output of f is a continuous variable this is called regression.

Both regression & classification are a form of function approximation.

Now also we could try to divide each input variables axis into small intervals, which splits \mathbb{R}^d into a grid of (hyper)-boxes. Each $\vec{x}^{(i)}$ in training set falls in a box and it has the label $y^{(i)}$.

Then, given a new example \vec{x}_{new} , we return $y = \{ \text{average of all training examples in that box} \}$

But if each axis has M divisions then there are M^d boxes but $M^d \gg \# \text{ training examples}$

So the data is very sparsely distributed.

So the naive approach won't work. That's one reason why dimension reduction is a valuable tool.

Fewer variables also reduces likelihood of overfitting.

In general, as model complexity grows so too does generalization error (caused by overfit).

Concretely, $P(\mathcal{C}_k) = \frac{\text{prior probability of character in } \mathcal{C}_k}$

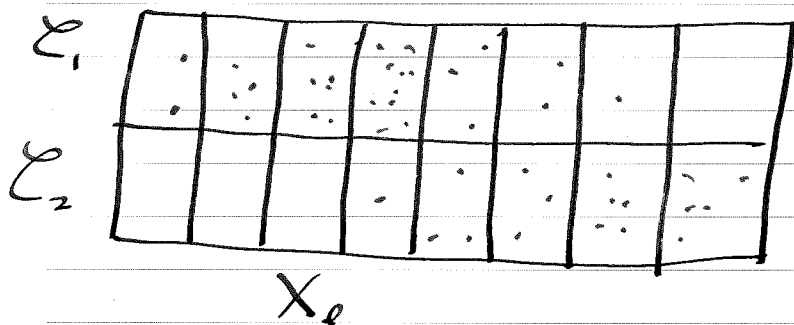
we can get these priors from the data themselves. For instance, if 'a' occurs 3 times more than 'b' then

$$P(\mathcal{C}_1) = \frac{3}{4} \quad P(\mathcal{C}_2) = \frac{1}{4}$$

if we had to classify a new character without being able to see it we should assign to the class with the higher probability.

Now suppose $\tilde{x}_i \in \{X_e\} \leftarrow$ discrete set of values, as used in histograms.

$P(\mathcal{C}_k, X_e) =$ probability the image has feature value X_e & belongs to class \mathcal{C}_k



$\therefore P(\mathcal{C}_k, X_e) =$ fraction of images in (\mathcal{C}_k, X_e) at cell

$P(\mathcal{C}_k) =$ fraction in each row

$P(X_e) =$ fraction in each column.

$P(X_e | \mathcal{C}_k) =$ probability of being in l 'th column given that you're in k 'th row

That's the function is now k fully in cell X_e .

$$P(\mathcal{C}_k | X_e) = \frac{P(X_e | \mathcal{C}_k) P(\mathcal{C}_k)}{P(X_e)}$$

by Bayes - as before since

$$P(\mathcal{C}_1 | X_e) + P(\mathcal{C}_2 | X_e) = 1$$

Then

$$\frac{P(X_e | \mathcal{C}_1) P(\mathcal{C}_1)}{P(X_e)} + \frac{P(X_e | \mathcal{C}_2) P(\mathcal{C}_2)}{P(X_e)} = 1$$

ie.

$$P(X_e) = P(X_e | \mathcal{C}_1) P(\mathcal{C}_1) + P(X_e | \mathcal{C}_2) P(\mathcal{C}_2)$$

ie.

$$P(\mathcal{C}_k | X_e) = \frac{P(X_e | \mathcal{C}_k) P(\mathcal{C}_k)}{P(X_e | \mathcal{C}_1) P(\mathcal{C}_1) + P(X_e | \mathcal{C}_2) P(\mathcal{C}_2)}$$

Notice everything on RHS is obtainable from the histogram of training data.

We should then say

$$f(\vec{x}) = \mathcal{C}_k \text{ where } P(\mathcal{C}_k | \vec{x}) > P(\mathcal{C}_j | \vec{x}) \quad \forall j \neq k$$

Each point of feature space then is assigned to one of N regions viz

$$P(x|\mathcal{C}_k) > P(x|\mathcal{C}_j) \quad j \neq k$$

ie.

$$f(x) = \underset{j \in \{1, \dots, N\}}{\operatorname{argmax}} P(x|\mathcal{C}_j)$$

This splits \mathbb{R}^d into c regions R_1, \dots, R_c

ie. $\mathbb{R}^d = \bigcup_{i=1}^c R_i$ Boundaries of which ∂R_i is
are decision boundaries or decision surfaces.

Now, consider again $N=2$

$$\begin{aligned} P(\text{error}) &= P(\tilde{x} \in R_2, \mathcal{C}_1) + P(x \in R_1, \mathcal{C}_2) \\ &= P(x \in R_2 | \mathcal{C}_1) P(\mathcal{C}_1) + P(x \in R_1 | \mathcal{C}_2) P(\mathcal{C}_2) \\ &= \int_{R_1} p(x|\mathcal{C}_1) P(\mathcal{C}_1) dx + \int_{R_2} p(x|\mathcal{C}_2) P(\mathcal{C}_2) dx \end{aligned}$$

if $p(x|\mathcal{C}_1) P(\mathcal{C}_1) > p(x|\mathcal{C}_2) P(\mathcal{C}_2)$ then we'd
classify it as C_1 , meaning R_1 should be R_2 should
be picked so that $\tilde{x} \in R_1$ since that integral is
smaller & contributes less error.

$y_k(\vec{x}) \stackrel{\text{def}}{=} P(\mathcal{E}_k | \vec{x})$ is called a discriminant function.

And our decisions are based on

$$f(\vec{x}) = \underset{j}{\operatorname{argmax}} \{y_j(\vec{x})\}$$

Of course regions R_i, R_j are contiguous

(i.e. $\partial R_i \cap \partial R_j \neq \emptyset$) then

$$\partial R_i \cap \partial R_j = \{ \vec{x} \in \mathbb{R}^d \mid y_i(\vec{x}) = y_j(\vec{x}) \}$$

We can characterize the 2-class classification problem in terms of

$$y(x) = y_1(\vec{x}) - y_2(\vec{x})$$

$$R_1 = \{ \vec{x} \mid y > 0 \} \quad R_2 = \{ \vec{x} \mid y < 0 \}$$

Discriminants can be found often without having to first do the probability estimations.

Generalization of this is with a loss matrix

$$L_{kj} = \left\{ \begin{array}{l} \text{penalty for assigning to } \mathcal{E}_j \\ \text{when should be } \mathcal{E}_k \end{array} \right\}$$

Then expected loss is

$$R_k = \sum_{j=1}^N L_{kj} \int_{R_j} p(x|C_k) dx$$

and the risk is $R = \sum_{k=1}^N R_k P(C_k)$

$$= \sum_{j=1}^N \int_{R_j} \sum_{k=1}^N L_{kj} p(x|C_k) P(C_k) dx$$

Risk is minimized when each integral is minimized
ie. when

$$f(\vec{x}) = \underset{j}{\text{argmin}} \sum_{k=1}^N L_{kj} p(x|C_k) P(C_k)$$

The prior decision rule used $L_{kj} = 1 - \delta_{kj}$
is indifferent to the kind of error. But sometimes
the type of misclassification matters more than that
(eg. medicine).