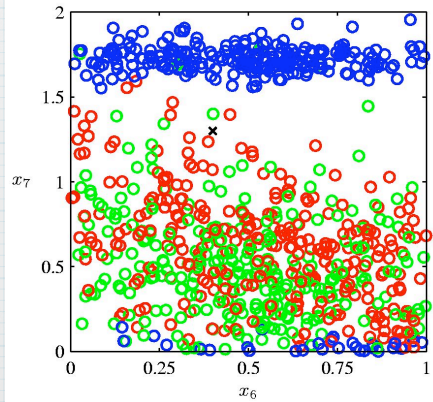


Curse of Dimensionality

- Real world applications deal with spaces with high dimensionality
- High dimensionality poses serious challenges for the design of pattern recognition techniques



Oil flow data in two dimensions

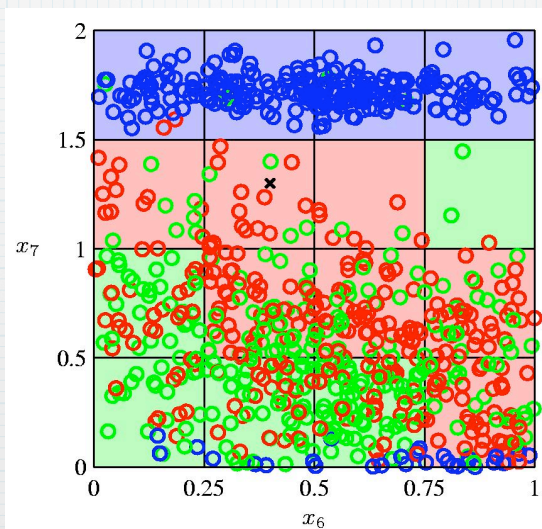
Red = homogeneous

Green = annular

Blue = laminar

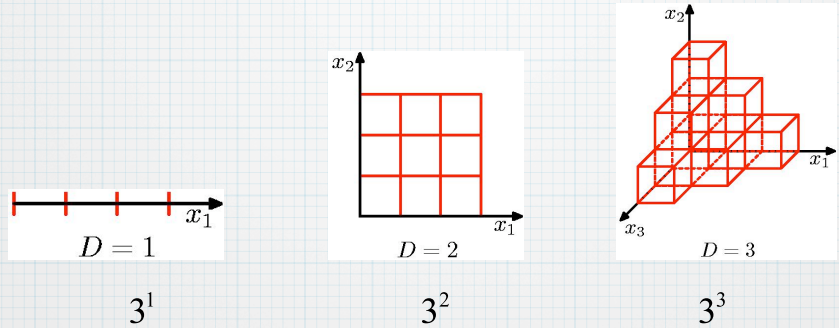
Curse of Dimensionality

- A simple classification approach



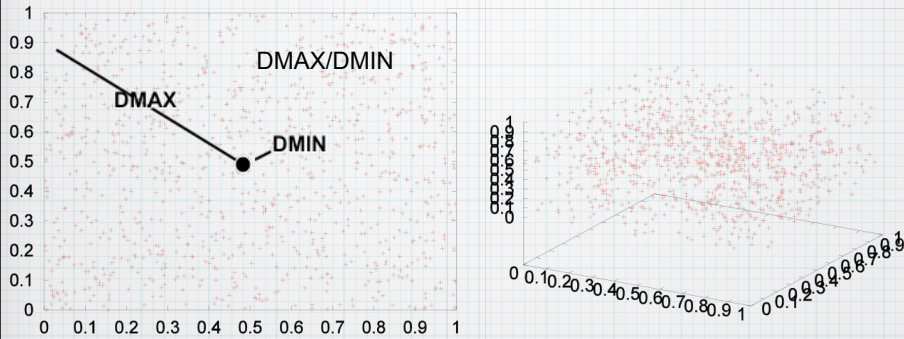
Curse of Dimensionality

➤ Going higher in dimensionality...



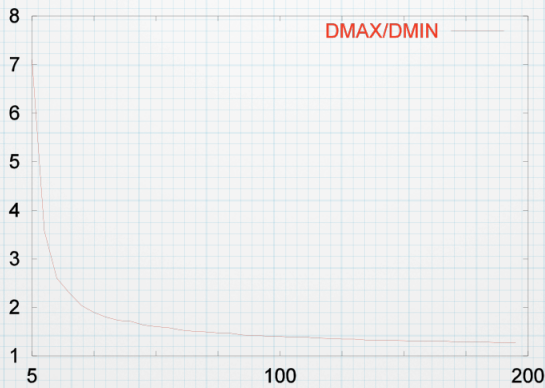
The number of regions of a regular grid grows exponentially with the dimensionality D of the space

Curse of Dimensionality



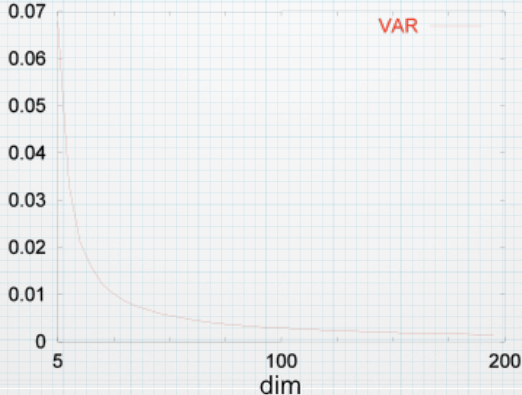
Sample of size $N=500$ uniformly distributed in $[0, 1]^q$

Curse of Dimensionality



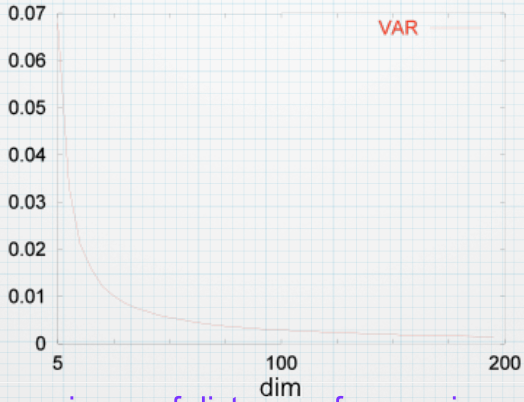
The distribution of the ratio **DMAX/DMIN** converges to 1 as the dimensionality increases

Curse of Dimensionality



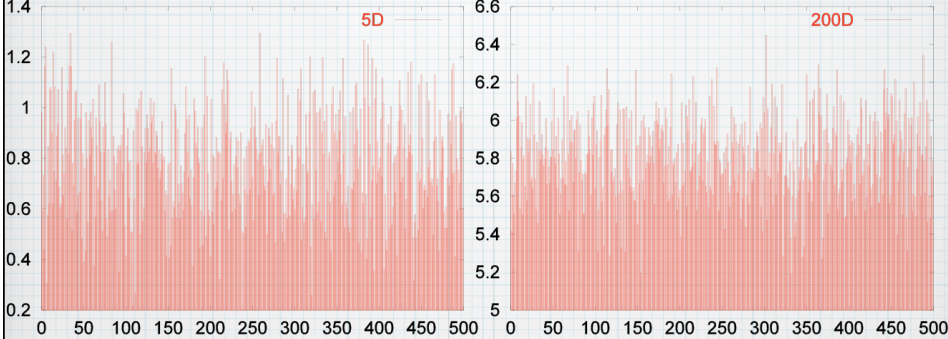
Variance of distances from a given point

Curse of Dimensionality

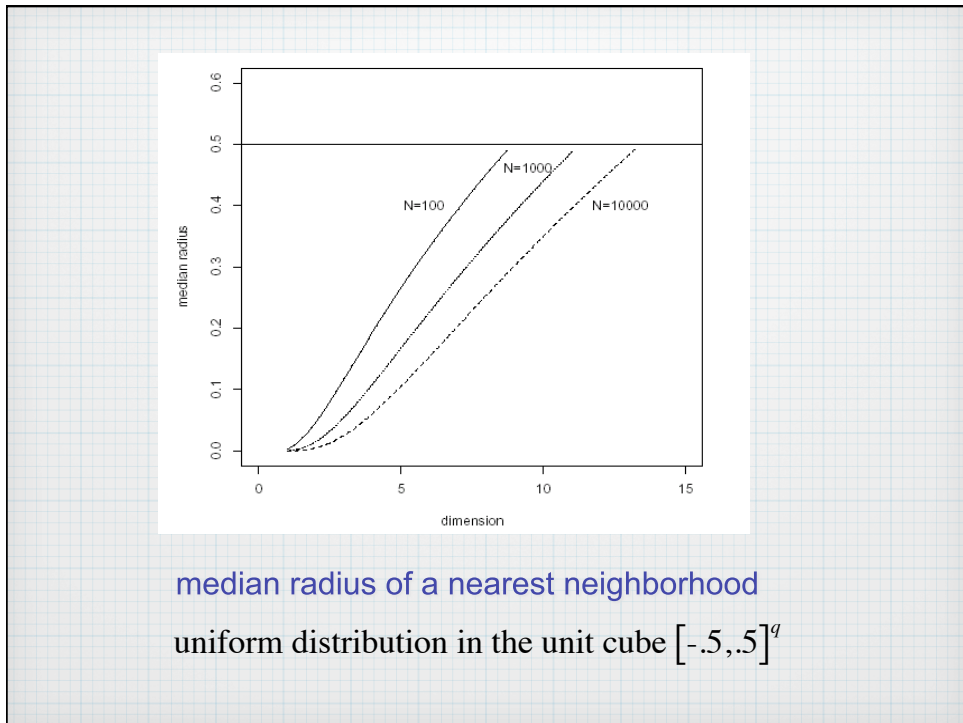
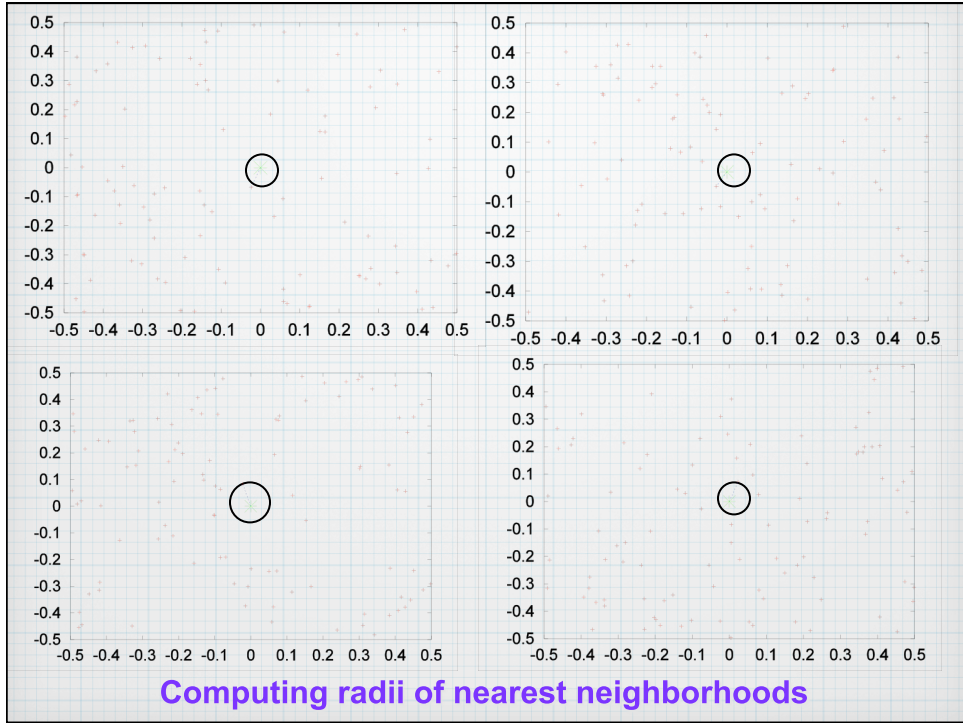


The variance of distances from a given point converges to 0 as the dimensionality increases

Curse of Dimensionality



Distance values from a given point
Values flatten out as dimensionality increases



Curse-of-Dimensionality

As dimensionality increases, the distance from the closest point increases faster

- Random sample of size $N \sim$ uniform distribution in the q -dimensional unit hypercube
- Diameter of a $K = 1$ neighborhood using Euclidean distance: $d(q,N) = O(N^{-1/q})$

| q | 4 | 4 | 6 | 6 | 10 | 10 | 20 | 20 | 20 |
|----------|------|------|------|------|------|-------|-------|--------|-----------|
| N | 100 | 1000 | 100 | 1000 | 1000 | 10000 | 10000 | 10^6 | 10^{10} |
| $d(q,N)$ | 0.42 | 0.23 | 0.71 | 0.48 | 0.91 | 0.72 | 1.51 | 1.20 | 0.76 |

Large $d(q,N) \Rightarrow$ Highly biased estimations

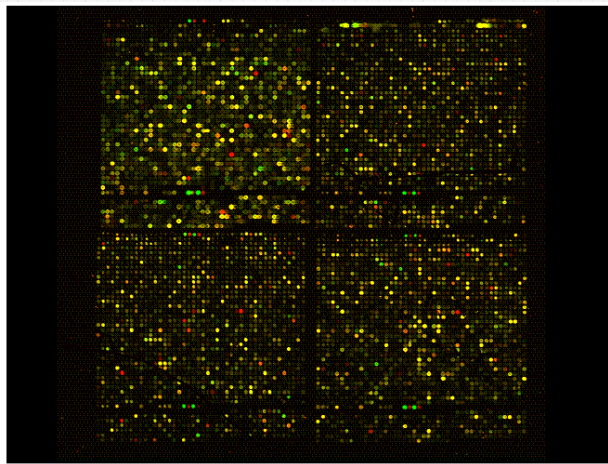
Curse-of-Dimensionality

- In high dimensional spaces data become extremely sparse and are far apart from each other
- **The curse of dimensionality affects any estimation problem with high dimensionality**

Curse-of-Dimensionality

- It is a serious problem in many real-world applications
- **Microarray data:** 3,000-4,000 genes;
- **Documents:** 10,000-20,000 words in dictionary;
- **Images, face recognition, etc.**

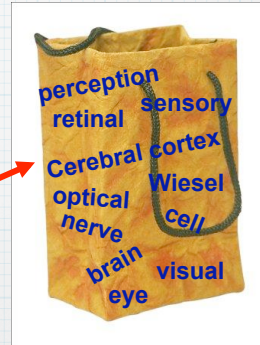
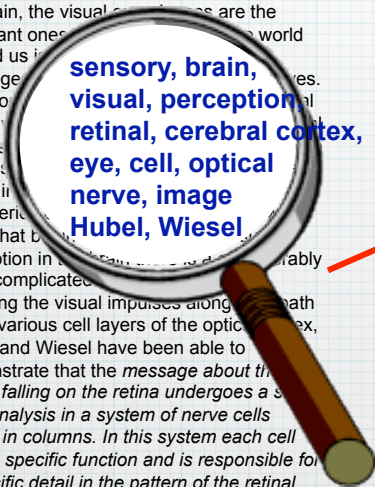
Microarray Data Analysis



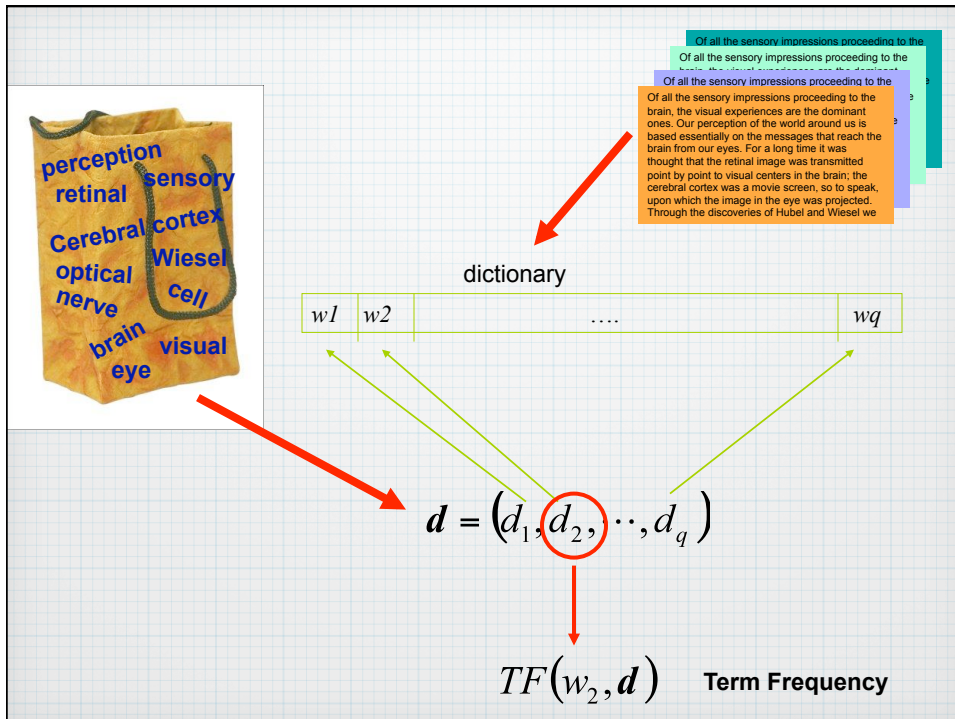
Problem: Which samples are most similar to each other, in terms of their expression profiles across genes

Document Classification

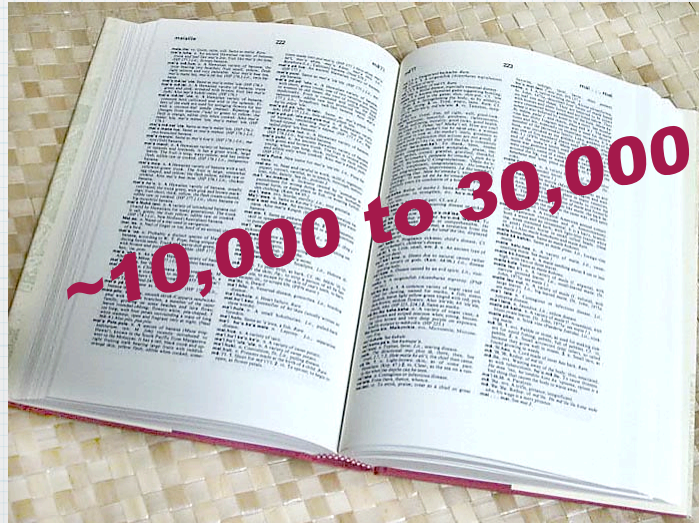
Of all the sensory impressions proceeding to the brain, the visual experiences are the dominant ones. Our perception of the world around us is based essentially on the messages that reach the brain from our eyes. For a long time it was thought that the retinal image was projected point by point to visual centers in the brain; the cerebral cortex was a movie screen, so to speak, upon which the image in the eye was projected. Through the discoveries of Hubel and Wiesel we know that this is not the case. Perception is a much more complicated process. It involves following the visual impulses along a path to the various cell layers of the optic cortex. Hubel and Wiesel have been able to demonstrate that the message about the image falling on the retina undergoes a *spatially selective analysis* in a system of nerve cells stored in columns. In this system each cell has its specific function and is responsible for a specific detail in the pattern of the retinal image.



Bag-of-words representation of a document



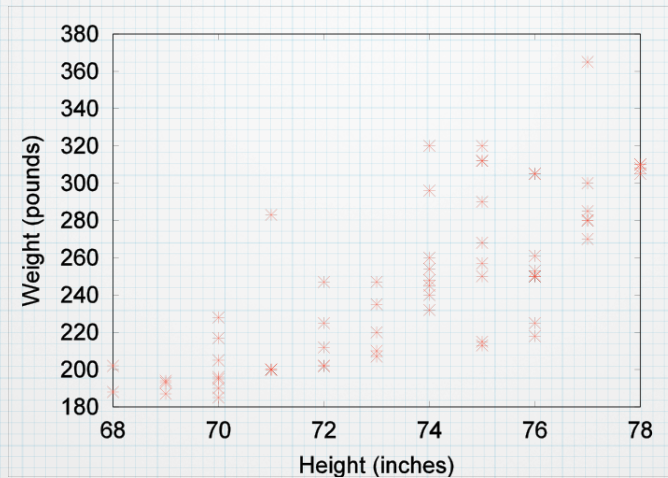
What's the size of the dictionary?



**How can we deal with
the curse of dimensionality?**

Curse-of-Dimensionality

- Effective techniques applicable to high dimensional spaces exist.
- The reasons are twofold:
 - ✓ Real data are often confined to regions of *lower dimensionality*
 - ✓ Real data typically exhibit *smoothness properties* (at least locally). Local interpolation techniques can be used to make predictions



$$\begin{bmatrix} 7.68 & 92.2 \\ 92.2 & 1912.5 \end{bmatrix}$$

$$\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \quad \boldsymbol{\mu} = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i$$

2 × 2 covariance matrix:

$$E[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T] =$$

$$E \left[\begin{pmatrix} x_1 - \mu_1 \\ x_2 - \mu_2 \end{pmatrix} (x_1 - \mu_1, x_2 - \mu_2) \right] =$$

$$E \begin{bmatrix} (x_1 - \mu_1)^2 & (x_1 - \mu_1)(x_2 - \mu_2) \\ (x_1 - \mu_1)(x_2 - \mu_2) & (x_2 - \mu_2)^2 \end{bmatrix} =$$

$$\frac{1}{N} \sum_{i=1}^N \begin{bmatrix} (x_1^i - \mu_1)^2 & (x_1^i - \mu_1)(x_2^i - \mu_2) \\ (x_1^i - \mu_1)(x_2^i - \mu_2) & (x_2^i - \mu_2)^2 \end{bmatrix}$$

$$\frac{1}{N} \sum_{i=1}^N \begin{bmatrix} (x_1^i - \mu_1)^2 & (x_1^i - \mu_1)(x_2^i - \mu_2) \\ (x_1^i - \mu_1)(x_2^i - \mu_2) & (x_2^i - \mu_2)^2 \end{bmatrix} =$$

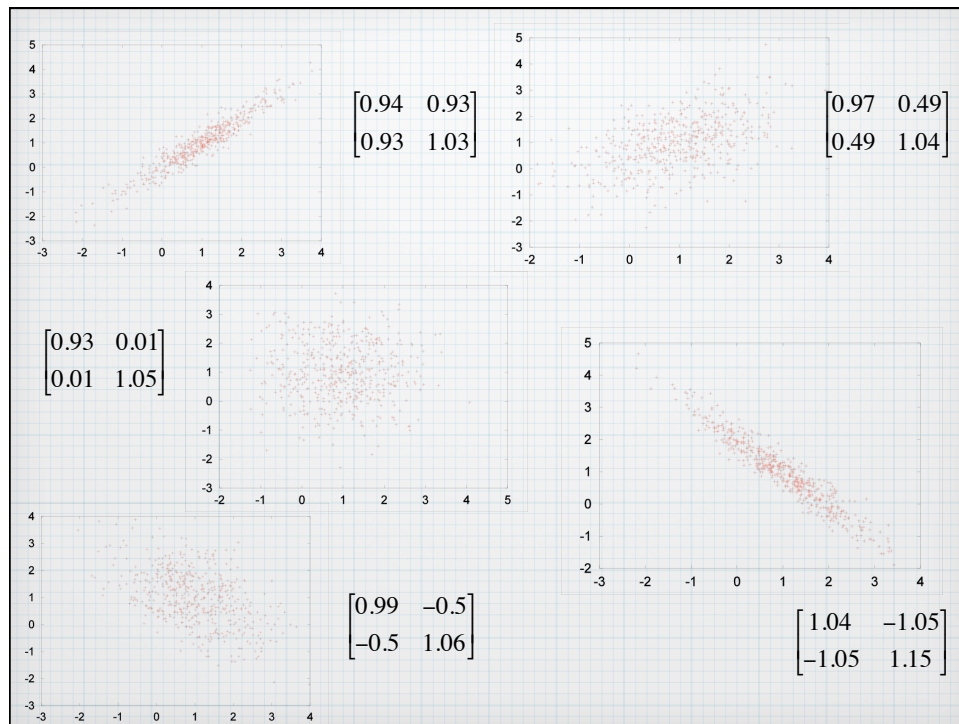
variance

covariance

$$\begin{bmatrix} \frac{1}{N} \sum_{i=1}^N (x_1^i - \mu_1)^2 & \frac{1}{N} \sum_{i=1}^N [(x_1^i - \mu_1)(x_2^i - \mu_2)] \\ \frac{1}{N} \sum_{i=1}^N [(x_1^i - \mu_1)(x_2^i - \mu_2)] & \frac{1}{N} \sum_{i=1}^N (x_2^i - \mu_2)^2 \end{bmatrix}$$

covariance

variance



Dimensionality Reduction

- Many dimensions are often interdependent (correlated);

We can:

- Reduce the dimensionality of problems;
- Transform interdependent coordinates into significant and independent ones;

Decision Theory

- *Decision theory*, when combined with *probability theory*, allows to make optimal decisions in situations involving uncertainty
- Training data: input vector \mathbf{x} , target vector \mathbf{t}
- Inference: joint probability distribution $p(\mathbf{x}, \mathbf{t})$
- Decision step: make optimal decision

Decision Theory

Classification example: medical diagnosis problem

- \mathbf{x} set of pixel intensities in an image
- Two classes:
 - $C_1 = 0$ absence of cancer
 - $C_2 = 1$ presence of cancer
- Inference step: estimate $p(\mathbf{x}, C_k)$
- Decision step: given \mathbf{x} predict C_k so that a measure of error is minimized according to the given probabilities

Decision Theory

How probabilities play a role in decision making?

- Decision step: given \mathbf{x} predict C_k

Thus, we are interested in $p(C_k | \mathbf{x})$

$$p(C_k | \mathbf{x}) = \frac{p(\mathbf{x} | C_k)p(C_k)}{p(\mathbf{x})}$$

Intuitively: we want to minimize the chance of assigning \mathbf{x} to the wrong class. Thus, choose the class that gives the higher posterior probability

Minimizing the misclassification rate

- Goal: Minimize the number of misclassifications

We need to find a rule that assigns each input vector to one of the possible classes C_k

Such rule divides the input space into regions R_k so that all points in R_k are assigned to C_k

Boundaries between regions are called **decision boundaries**

Minimizing the misclassification rate

- Goal: Minimize the number of misclassifications

$$\begin{aligned} p(\text{mistake}) &= p(x \in R_1, C_2) + p(x \in R_2, C_1) \\ &= \int_{R_1} p(x, C_2) dx + \int_{R_2} p(x, C_1) dx \end{aligned}$$

- Assign x to the class that gives the smaller value of the integrand:

- Choose C_1 if $p(x, C_1) > p(x, C_2)$
- Choose C_2 if $p(x, C_2) > p(x, C_1)$

Minimizing the misclassification rate

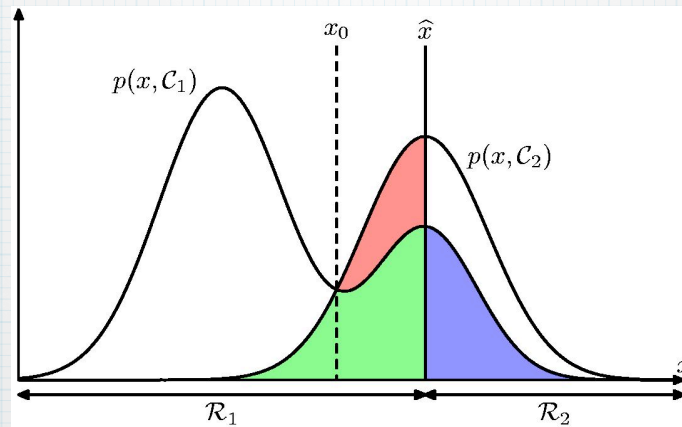
- Choose C_1 if $p(x, C_1) > p(x, C_2)$
- Choose C_2 if $p(x, C_2) > p(x, C_1)$

$$p(x, C_k) = p(C_k | x)p(x)$$

Thus:

- Choose C_1 if $p(C_1 | x) > p(C_2 | x)$
- Choose C_2 if $p(C_2 | x) > p(C_1 | x)$

Minimizing the misclassification rate



Optimal decision boundary: $\hat{x} = x_0$

Minimizing the misclassification rate

General case of K classes:

$$p(\text{correct}) = \sum_{k=1}^K p(\mathbf{x} \in R_k, C_k) = \sum_{k=1}^K \int_{R_k} p(\mathbf{x}, C_k) d\mathbf{x}$$

Thus:

Choose C_k that gives the largest $p(C_k | \mathbf{x})$

Minimizing the expected loss

- Some mistakes are more costly than others.
- **Loss function (cost function)**: overall measure of loss incurred in taking any of the available decisions

L_{kj} : loss incurred when we assign \mathbf{x} to class C_j and the true class is C_k

| | cancer | normal |
|--------|--------|--------|
| cancer | 0 | 1000 |
| normal | 1 | 0 |

The optimal solution is the one that minimizes the loss function

Minimizing the expected loss

- The loss function depends on the true class, which is unknown.
- The uncertainty of the true class is expressed through the joint probability $p(\mathbf{x}, C_k)$
- We minimize the expected loss:

$$E[L] = \sum_k \sum_j \int_{R_j} L_{kj} p(\mathbf{x}, C_k) d\mathbf{x}$$

- For each \mathbf{x} we should minimize

$$\sum_k L_{kj} p(\mathbf{x}, C_k) = \sum_k L_{kj} p(C_k | \mathbf{x}) p(\mathbf{x})$$

Minimizing the expected loss

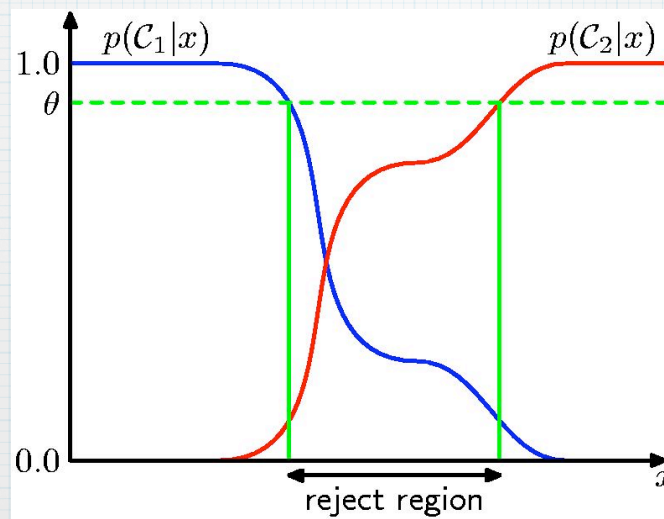
- For each x we should minimize

$$\sum_k L_{kj} p(\mathbf{x}, C_k) = \sum_k L_{kj} p(C_k | \mathbf{x}) p(\mathbf{x})$$

- Thus, to minimize the expected loss: Assign each x to the class j that minimizes

$$\sum_k L_{kj} p(C_k | \mathbf{x})$$

The Reject Option



Inference and Decision

- **Inference stage**: use the training data to learn a model for $p(C_k | \mathbf{x})$

- **Decision stage**: use the given posterior probabilities to make optimal class assignments

Generative Methods

- Solve the inference problem of estimating the class-conditional densities $p(\mathbf{x} | C_k)$ for each class C_k
- Infer the prior class probabilities $p(C_k)$
- Use Bayes' theorem to find the class posterior probabilities:

$$p(C_k | \mathbf{x}) = \frac{p(\mathbf{x} | C_k)p(C_k)}{p(\mathbf{x})}$$
 where

$$p(\mathbf{x}) = \sum_k p(\mathbf{x} | C_k)p(C_k)$$
- Use decision theory to determine class membership for each new input \mathbf{x}

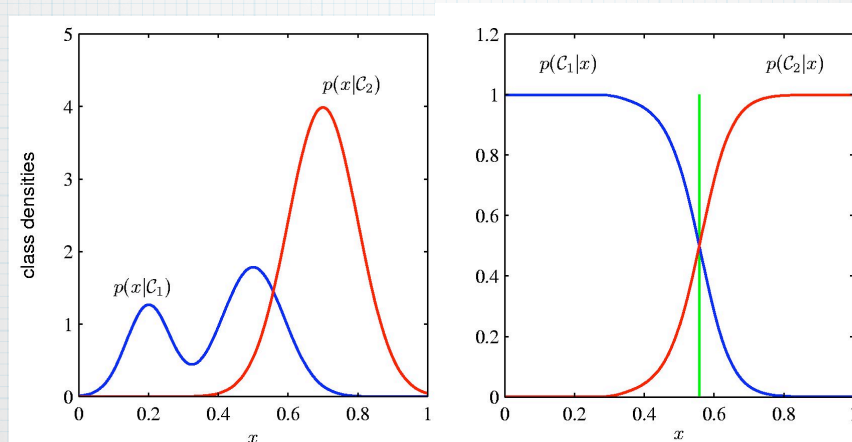
Discriminative Methods

- Solve directly the inference problem of estimating the class posterior probabilities $p(C_k | \mathbf{x})$
- Use decision theory to determine class membership for each new input \mathbf{x}

Discriminant Functions

- Find a function $f(\mathbf{x})$ which maps each input directly onto a class label. Probabilities play no role here.
- Use decision theory to determine class membership for each new input \mathbf{x}

Example



Linear Models for Classification

- **Classification:** Given an input vector x , assign it to one of K classes C_k where $k = 1, \dots, K$
- The input space is divided in **decision regions** whose boundaries are called **decision boundaries** or **decision surfaces**
- **Linear models:** decision surfaces are linear functions of the input vector x . They are defined by $(D - 1)$ -dimensional hyperplanes within the D -dimensional input space

Linear Models for Classification

- For regression: $y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$
- For classification, we want to predict class labels, or more generally class posterior probabilities.
- We transform the linear function of \mathbf{w} using a nonlinear function $f()$ so that

$$f(\mathbf{w}^T \mathbf{x} + w_0)$$

Generalized Linear Models

Linear Discriminant Functions

Two classes:

$$y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$$

if $y(\mathbf{x}) \geq 0$ *assign* \mathbf{x} *to* C_1

otherwise *assign* \mathbf{x} *to* C_2

Decision boundary: $y(\mathbf{x}) = 0$

Linear Discriminant Functions

Geometrical properties:

Decision boundary: $y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 = 0$

Let $\mathbf{x}_1, \mathbf{x}_2$ be two points which lie on the decision boundary

$$y(\mathbf{x}_1) = \mathbf{w}^T \mathbf{x}_1 + w_0 = 0, y(\mathbf{x}_2) = \mathbf{w}^T \mathbf{x}_2 + w_0 = 0$$

$$\Rightarrow \mathbf{w}^T (\mathbf{x}_1 - \mathbf{x}_2) = 0$$

\mathbf{w} represents the orthogonal direction
to the decision boundary

Geometrical properties (con't)

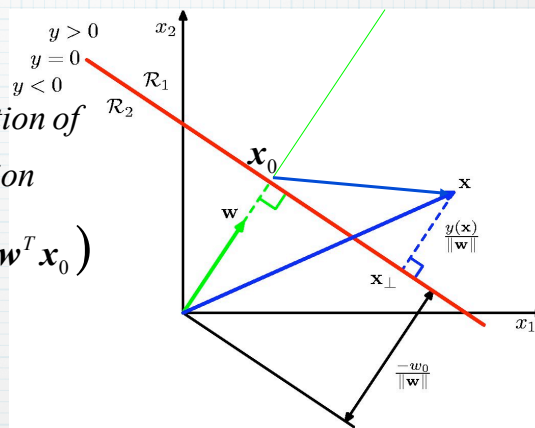
$$\mathbf{w}^{*T} = \frac{\mathbf{w}^T}{\|\mathbf{w}\|}$$

$\mathbf{w}^{*T} (\mathbf{x} - \mathbf{x}_0)$ is the projection of
 $(\mathbf{x} - \mathbf{x}_0)$ onto the \mathbf{w}^* direction

$$\frac{\mathbf{w}^T}{\|\mathbf{w}\|} (\mathbf{x} - \mathbf{x}_0) = \frac{1}{\|\mathbf{w}\|} (\mathbf{w}^T \mathbf{x} - \mathbf{w}^T \mathbf{x}_0)$$

$$= \frac{1}{\|\mathbf{w}\|} (\mathbf{w}^T \mathbf{x} + w_0) = \frac{y(\mathbf{x})}{\|\mathbf{w}\|}$$

when $\mathbf{x} = \mathbf{0}$, $\frac{y(\mathbf{x})}{\|\mathbf{w}\|} = \frac{w_0}{\|\mathbf{w}\|}$

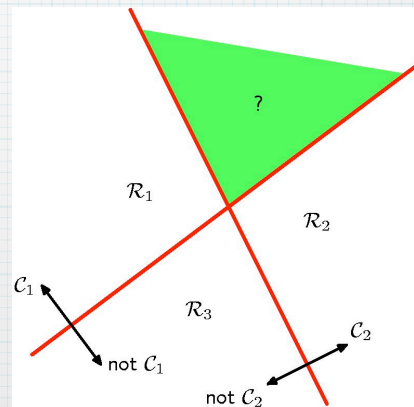


Signed orthogonal distance of
the origin from the decision
surface

Linear Discriminant Functions

Multiple classes

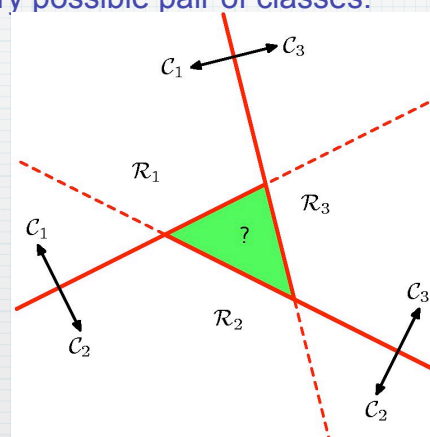
one-versus-the-rest: $K-1$ classifiers each of which solves a two-class problem of separating points of C_k from points not in that class



Linear Discriminant Functions

Multiple classes

one-versus-one: $K(K-1)/2$ binary discriminant functions, one for every possible pair of classes.



Linear Discriminant Functions

Multiple classes

Solution: consider a single K -class discriminant comprising K linear functions of the form

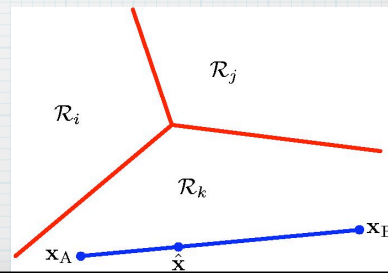
$$y_k(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x} + w_{k0}$$

Assign a point \mathbf{x} to class C_k if $y_k(\mathbf{x}) > y_j(\mathbf{x}) \forall j \neq k$

The decision boundary between class C_k and class C_j is given by

$$y_k(\mathbf{x}) = y_j(\mathbf{x})$$

$$\Rightarrow (\mathbf{w}_k - \mathbf{w}_j)^T \mathbf{x} + (w_{k0} - w_{j0}) = 0$$



Linear Discriminant Functions

Two approaches:

- Fisher's linear discriminant
- Perceptron algorithm

Fisher's Linear Discriminant

One way to view a linear classification model is in terms of *dimensionality reduction*.

Two class case:

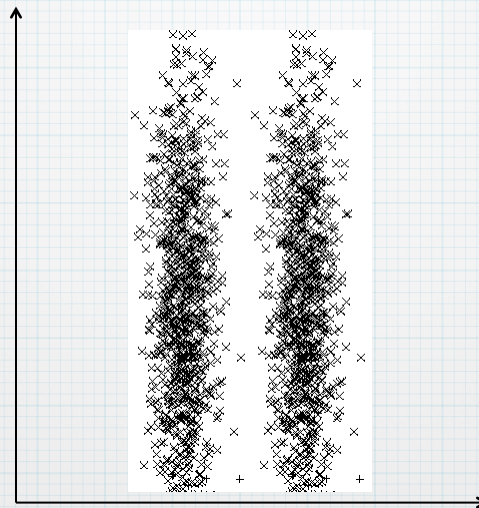
Suppose we project \mathbf{x} onto one dimension:

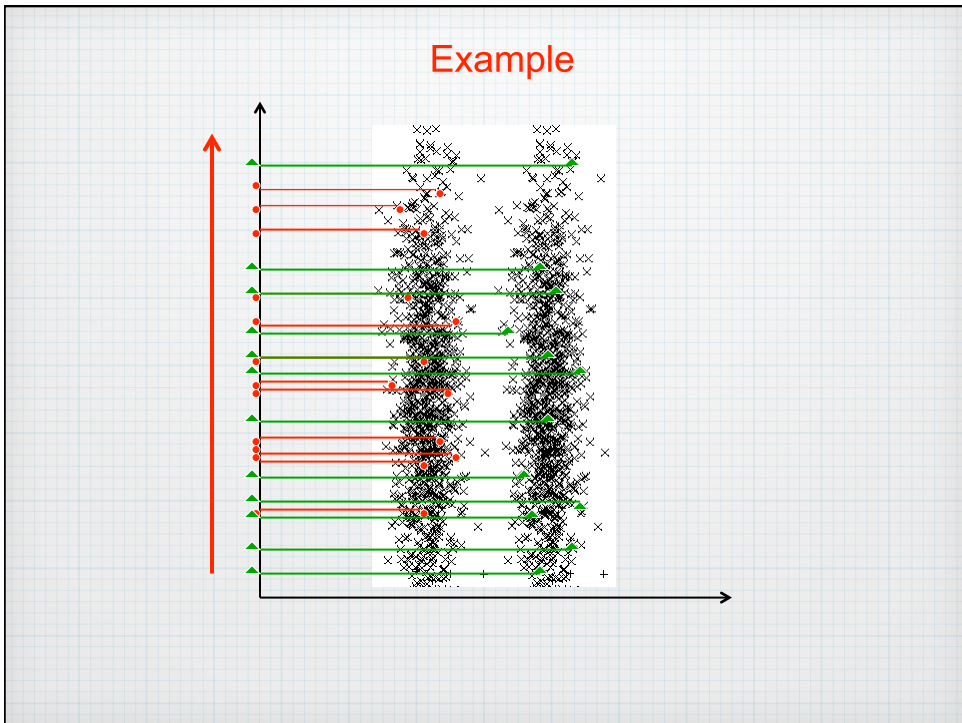
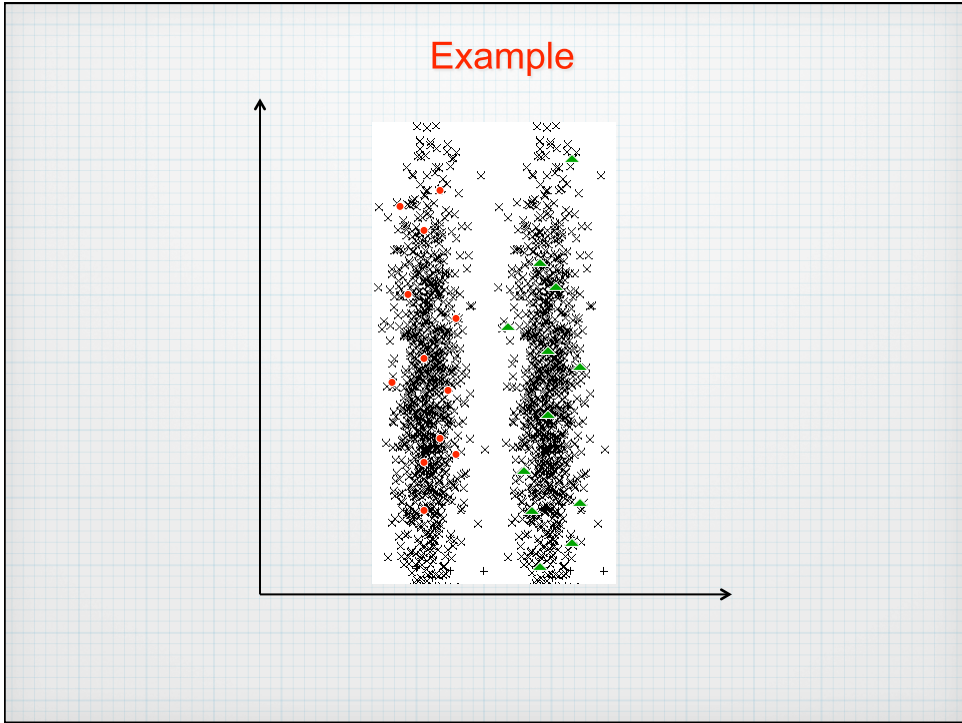
$$y = \mathbf{w}^T \mathbf{x}$$

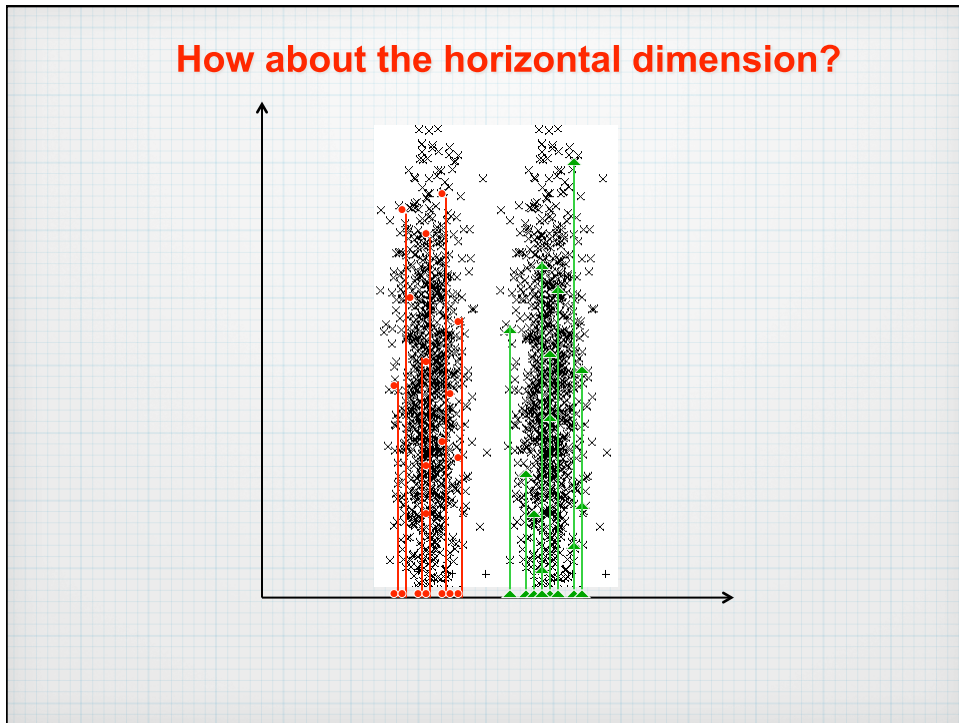
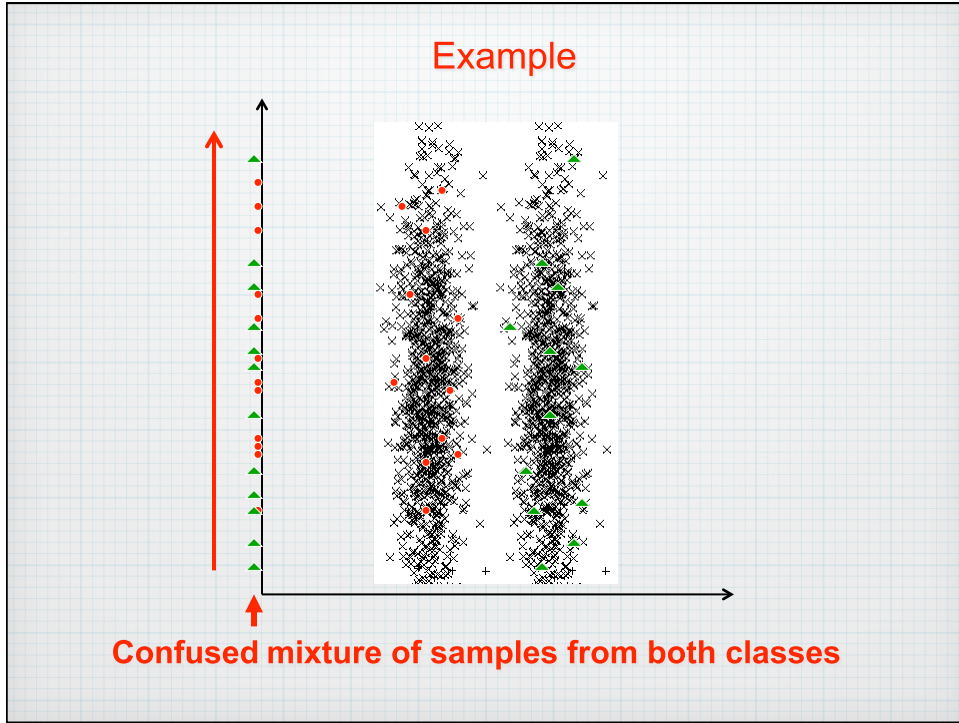
Set a threshold t

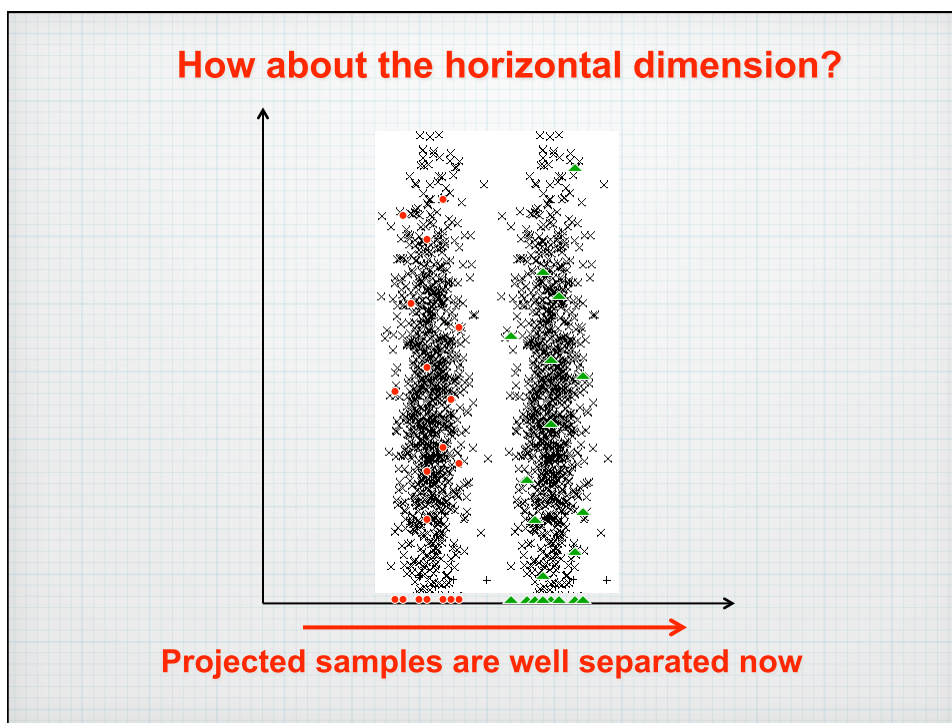
if $y \geq t$ assign \mathbf{x} to C_1
otherwise assign \mathbf{x} to C_2

Example









Fisher's Linear Discriminant

- Find an orientation along which the projected samples are well separated;
- This is exactly the goal of **linear discriminant analysis (LDA)**;
- In other words: we are after the linear projection that best separates the data, i.e. best **discriminates** data of different classes.

How can we find such discriminant direction?

LDA

$$\{(\mathbf{x}_n, C_i)\}_{i=1}^N \quad \mathbf{x}_n \in \mathbb{R}^q \quad C_i \in \{C_1, C_2\}$$

- N_1 samples of class C_1
- N_2 samples of class C_2
- Consider $\mathbf{w} \in \mathbb{R}^q$ with $\|\mathbf{w}\| = 1$
- Then: $\mathbf{w}^T \mathbf{x}$ is the projection of \mathbf{x} along the direction of \mathbf{w}
- We want the projections $\mathbf{w}^T \mathbf{x}$ where $\mathbf{x} \in C_1$ separated from the projections $\mathbf{w}^T \mathbf{x}$ where $\mathbf{x} \in C_2$

LDA

- A measure of the separation between the projected points is the difference of the sample means:

$$\mathbf{m}_i = \frac{1}{N_i} \sum_{\mathbf{x} \in C_i} \mathbf{x} \quad \text{Sample mean of class } C_i$$

$$m_i = \frac{1}{N_i} \sum_{\mathbf{x} \in C_i} \mathbf{w}^T \mathbf{x} = \mathbf{w}^T \mathbf{m}_i \quad \text{Sample mean for the projected points}$$

$$\Rightarrow \quad |m_1 - m_2| = |\mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2)|$$

We wish to make the above difference as large as we can. In addition...

LDA

- To obtain good separation of the projected data we really want the difference between the means to be large relative to some measure of the standard deviation of each class:

$$s_i^2 = \sum_{x \in C_i} (w^T x - m_i)^2$$

Scatter for the projected samples of class C_i

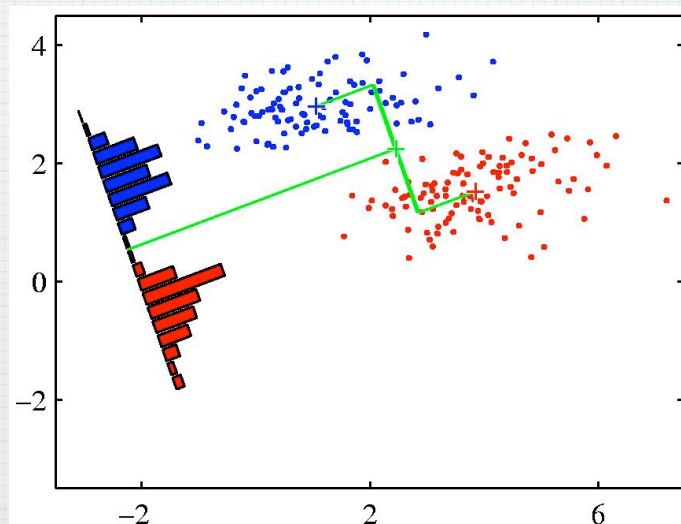
$$s_1^2 + s_2^2$$

Total **within-class scatter** of the projected samples

$$\arg \max_w \frac{|m_1 - m_2|^2}{s_1^2 + s_2^2}$$

Fisher linear discriminant analysis

LDA



LDA

$$J(\mathbf{w}) = \frac{|m_1 - m_2|^2}{s_1^2 + s_2^2}$$

To obtain $J(\mathbf{w})$ as an explicit function of \mathbf{w} we define the following matrices :

$$S_i = \sum_{x \in C_i} (\mathbf{x} - \mathbf{m}_i)(\mathbf{x} - \mathbf{m}_i)^T$$

$$S_W = S_1 + S_2 \quad \text{Within-class scatter matrix}$$

Then:

$$\begin{aligned} s_i^2 &= \sum_{x \in C_i} (\mathbf{w}^T \mathbf{x} - m_i)^2 = \sum_{x \in C_i} (\mathbf{w}^T \mathbf{x} - \mathbf{w}^T \mathbf{m}_i)^2 \\ &= \sum_{x \in C_i} \mathbf{w}^T (\mathbf{x} - \mathbf{m}_i)(\mathbf{x} - \mathbf{m}_i)^T \mathbf{w} = \mathbf{w}^T S_i \mathbf{w} \end{aligned}$$

LDA

$$\text{So: } s_1^2 = \mathbf{w}^T S_1 \mathbf{w} \quad \text{and} \quad s_2^2 = \mathbf{w}^T S_2 \mathbf{w}$$

$$\begin{aligned} \text{Thus: } s_1^2 + s_2^2 &= \mathbf{w}^T S_1 \mathbf{w} + \mathbf{w}^T S_2 \mathbf{w} = \\ &= \mathbf{w}^T (S_1 + S_2) \mathbf{w} = \mathbf{w}^T S_W \mathbf{w} \end{aligned}$$

Similarly :

$$\begin{aligned} (m_1 - m_2)^2 &= (\mathbf{w}^T \mathbf{m}_1 - \mathbf{w}^T \mathbf{m}_2)^2 = \\ &= \mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2)(\mathbf{m}_1 - \mathbf{m}_2)^T \mathbf{w} = \\ &= \mathbf{w}^T S_B \mathbf{w} \end{aligned}$$

$$\text{where } S_B = (\mathbf{m}_1 - \mathbf{m}_2)(\mathbf{m}_1 - \mathbf{m}_2)^T \quad \text{Between-class scatter matrix}$$

LDA

We have obtained :

$$s_1^2 + s_2^2 = \mathbf{w}^T S_W \mathbf{w}$$

$$(m_1 - m_2)^2 = \mathbf{w}^T S_B \mathbf{w}$$

$$\Rightarrow J(\mathbf{w}) = \frac{|m_1 - m_2|^2}{s_1^2 + s_2^2} = \frac{\mathbf{w}^T S_B \mathbf{w}}{\mathbf{w}^T S_W \mathbf{w}}$$

$$\arg \max_{\mathbf{w}} \frac{\mathbf{w}^T S_B \mathbf{w}}{\mathbf{w}^T S_W \mathbf{w}}$$

LDA

$$J(\mathbf{w}) = \frac{\mathbf{w}^T S_B \mathbf{w}}{\mathbf{w}^T S_W \mathbf{w}}$$

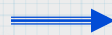
$$J(\mathbf{w}) \text{ is maximized when } (\mathbf{w}^T S_B \mathbf{w}) S_W \mathbf{w} = (\mathbf{w}^T S_W \mathbf{w}) S_B \mathbf{w}$$

We observe that:

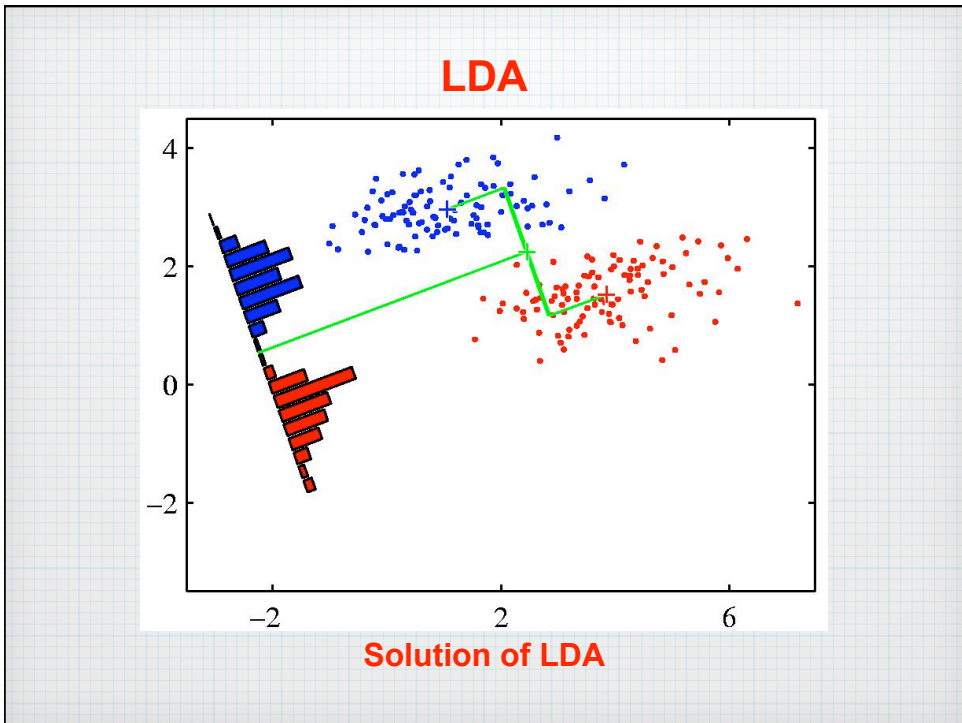
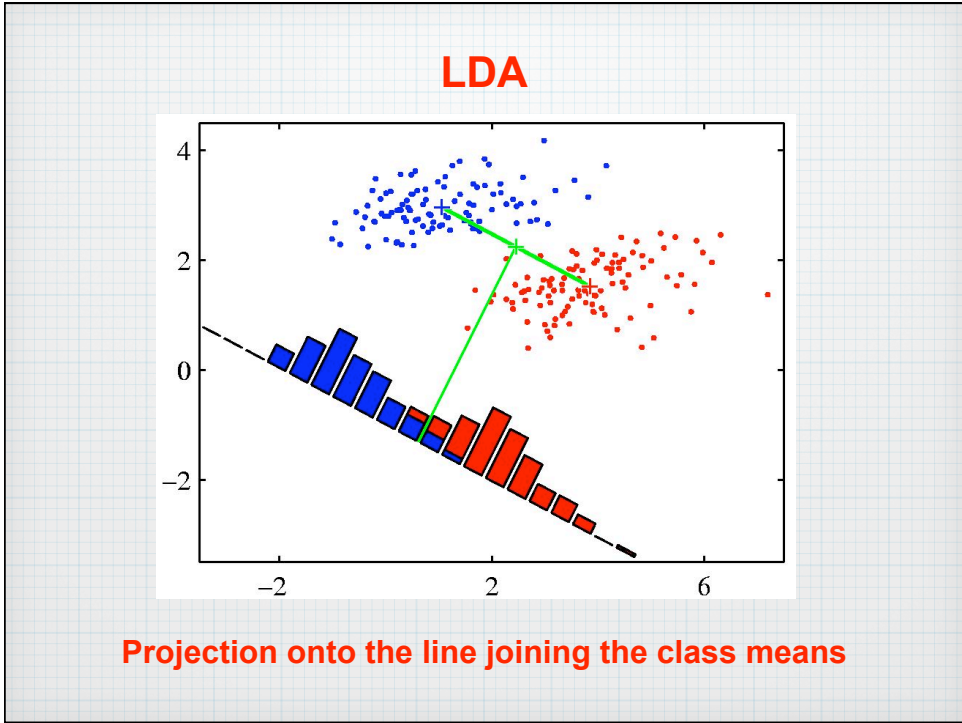
$$S_B \mathbf{w} = (m_1 - m_2)(m_1 - m_2)^T \mathbf{w}$$

← scalar →

Always in the direction of
← (m₁ - m₂) →



$$\mathbf{w} = S_W^{-1} (m_1 - m_2)$$



LDA

$$\mathbf{w} = S_W^{-1}(\mathbf{m}_1 - \mathbf{m}_2)$$

- Gives the linear function with the maximum ratio of between-class scatter to within-class scatter.
- The problem, e.g. classification, has been reduced to a q -dimensional problem to a more manageable one-dimensional problem.
- Optimal for multivariate normal class conditional densities.

LDA

- The analysis can be extended to multiple classes.
- LDA is a **linear** technique for dimensionality reduction: it projects the data along directions that can be expressed as **linear combination** of the input features.
- The “appropriate” transformation depends on the data and on the **task** we want to perform on the data. Note that LDA uses class labels.
- Non-linear extensions of LDA exist (e.g., generalized LDA).