

High Dimensional

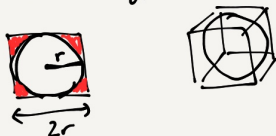
The quality of BDR depends on the CCD estimates.

How does it work when X is high-dimensional?

"High dimensional spaces are weird!"
(do not trust your intuition.)

Examples:

(1) consider a hypercube & an inscribed hypersphere in \mathbb{R}^d .



volume of hypersphere: $V_d(r) = \frac{\pi^{d/2} r^d}{\Gamma(\frac{d}{2} + 1)}$

\leftarrow Gamma function

$$\Gamma(n) = \int_0^\infty e^{-x} x^{n-1} dx$$

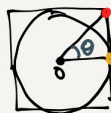
$$\Gamma(n+1) = n!$$

volume of hypercube: $(2r)^d$

let $f_d = \frac{\text{volume sphere}}{\text{volume cube}} = \frac{\pi^{d/2}}{2^d \Gamma(\frac{d}{2} + 1)}$

d	1	2	3	$\rightarrow \infty$
f_d	1	0.785	0.524	$\rightarrow 0$

as d increases, the volume of corners increases.
(of cube)



$$c = [r, r, r, \dots, r]$$

$$p = [r, 0, 0, \dots, 0]$$

$$\|c\|^2 = dr^2$$

$$\|p\|^2 = r^2$$

$$\cos \theta = \frac{c^T p}{\|c\| \|p\|} = \frac{r^2}{\sqrt{d} r^2} = \frac{1}{\sqrt{d}}$$

as d increases then $c \perp p$!
(corner is orthogonal to the axis)

$d=1$



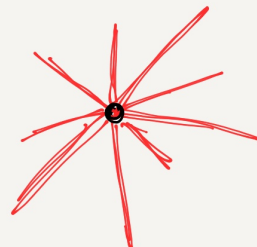
$d=2$



$d=3$

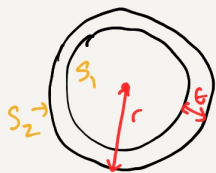


$d = \text{large}$



Example 2

consider a hypersphere shell of thickness G .



$$V_{\text{shell}} = V(S_2) - V(S_1) \\ = \left(1 - \frac{V(S_1)}{V(S_2)}\right) V(S_2)$$

$$\frac{V(S_1)}{V(S_2)} = \frac{\frac{(r-G)^{d/2} \pi^{d/2}}{\Gamma(\frac{d}{2}+1)}}{\frac{r^{d/2} \pi^{d/2}}{\Gamma(\frac{d}{2}+1)}} = \left(1 - \frac{G}{r}\right)^d$$

< 1

suppose $0 < G < r$,
as d increases, then $\frac{V(S_1)}{V(S_2)} \rightarrow 0$

$\Rightarrow V_{\text{shell}} \rightarrow V(S_2)$ as d increases

"All the volume is in the shell of the hypersphere"

Example 3

 high-dim Gaussian

let $X \sim N(0, \sigma^2 I)$

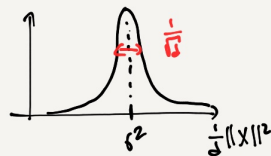
i.e. $x_i \sim N(0, \sigma^2)$ iid c.v.

$$\text{Then, } E[\|X\|^2] = E[x_1^2 + x_2^2 + \dots + x_d^2] = d\sigma^2$$

$$E\left[\frac{1}{d}\|X\|^2\right] = \sigma^2$$

Note: $\|X\|^2$ is a sum of iid c.v., thus by the central limit theorem it is concentrated around the mean as $d \rightarrow \infty$.

$$\frac{1}{d}\|X\|^2 \sim N(\sigma^2, \frac{1}{d})$$

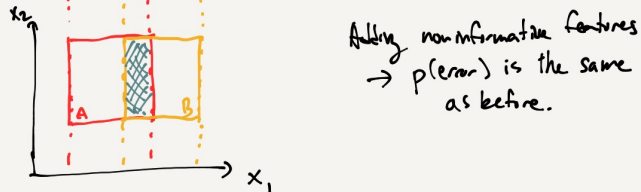
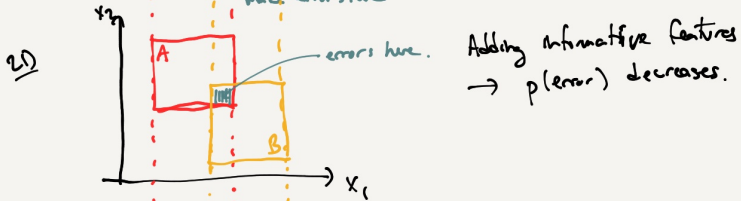
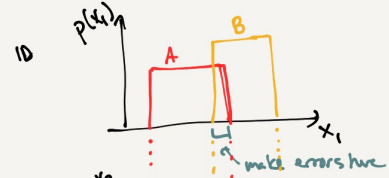


In high-dim, a Gaussian is essentially a shell of radius $\sqrt{d}\sigma$. Most of the density is in the shell.

(max density is still the mean)

Curse of Dimensionality

In theory, adding new features will not increase $P(\text{error})$



In practice, for BDR, error increases as feat dim increases.

The problem: Quality of the CCD estimates.

Density estimates in high-dim require more training samples.

Roughly, desired training set size = $O(e^p)$, $p = \#$ of parameters

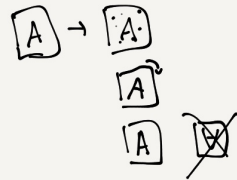
Solution:

- 1) Reduce $\#$ of parameters (complexity of model)
(e.g. full cov \rightarrow diag cov)
- \rightarrow 2) Reduce $\#$ of features (dimensionality reduction)
 \rightarrow implicitly reduce $\#$ of parameters.

- 3) Create more data

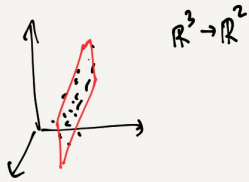
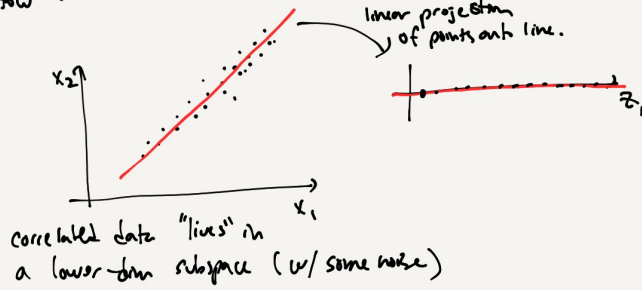
a) Bayesian estimation (virtual samples)

b) data augmentation



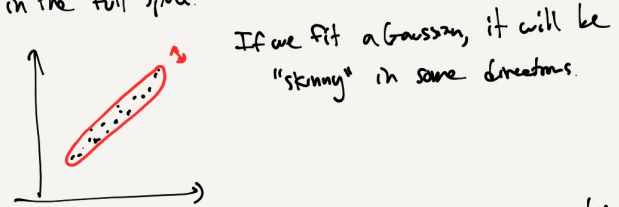
Linear Dimensionality Reduction

- Summarize correlated features w/ fewer features
- How to find these correlations?



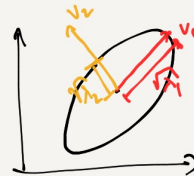
Principal Component Analysis (PCA)

Idea: if the data lives in a subspace, then it will look flat in the full space.



let (v_i, λ_i) be an eigenpair of covariance matrix Σ
 $\Sigma = V \Lambda V^T$, $V = [v_1, \dots, v_d]$, $\Lambda = \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_d \end{bmatrix}$

- each v_i defines an axis of ellipse
- each λ_i defines the width on that axis



Hence, the eigenvalues of Σ tell us which directions the data is flat.

\Rightarrow select axis v_i w/ larger eigenvalues as "principal components".

PCA: Given the dataset $\{x_1, \dots, x_n\}$ of dim k

- training**
- 1) Calculate Gaussian

$$\mu = \frac{1}{n} \sum_{i=1}^n x_i, \quad \Sigma = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)(x_i - \mu)^T$$
 - 2) eigen decomp of Σ : $\Sigma = V \Lambda V^T$
 - 3) order the eigenvalues: $\lambda_1 > \lambda_2 > \lambda_3 > \dots > \lambda_d > 0$
 - 4) select the top- k eigenvectors: $\Phi = [v_1, \dots, v_k]$
- dim. reduction**
- 5) project new point x onto Φ :

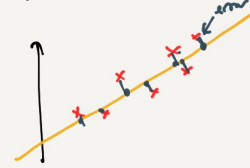
$$z = \Phi^T (x - \mu) \leftarrow \text{PCA coefficients. (new feature vector)}$$

(use BDR or other classifiers)

Notes:

The selection of Φ w/ $\Phi^T \Phi = I$ also:

- (1) maximizes the variance of the projected data in z . PS7-3
- (2) minimizes the reconstruction error of training data. PS7-2



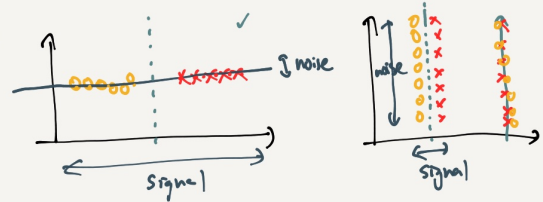
(3) can be implemented efficiently w/ SVD PS7-4 (tutorial)

(4) select k ? 1) pick k that works in the downstream task (classification)

2) pick k to preserve variance of data

$$p\% = \frac{\sum_{i=1}^k \lambda_i}{\sum_{i=1}^d \lambda_i} \quad (95\%)$$

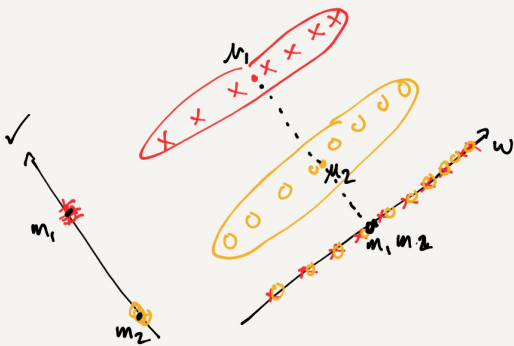
(5)



Assumption that signal variance is larger than the noise variance.

(6) PCA optimal for representation (but not necessarily for classification)

Fisher's Linear Discriminant (FLD) (Linear Discriminant Analysis (LDA))



Find the projection that best separates the classes.
 $z = w^T x$

Class statistics
class mean

original space
 $\mu_j = \frac{1}{n_j} \sum_{x_i \in C_j} x_i$

1-d space
 $m_j = w^T \mu_j$

class scatter

$S_j = \sum_{x_i \in C_j} (x_i - \mu_j)(x_i - \mu_j)^T$

$S_j = w^T S_j w$

IDEA: maximize the distance b/w projected means
 $(m_1 - m_2)^2 = (w^T (\mu_1 - \mu_2))^2$

problem: w is unconstrained \rightarrow need normalization

Fisher's Idea

$$w^* = \underset{w}{\operatorname{argmax}} \frac{\overbrace{(\mu_1 - \mu_2)^2}^{\text{between-class scatter}}}{\underbrace{S_1 + S_2}_{\text{within class scatter}}} = \underset{w}{\operatorname{argmax}} \frac{w^T S_B w}{w^T S_W w}$$

$$S_B = (\mu_1 - \mu_2)(\mu_1 - \mu_2)^T$$

$$S_W = S_1 + S_2$$

\downarrow generalized eigenvalue problem

$$w^* = (S_1 + S_2)^{-1} (\mu_1 - \mu_2)$$

"Fisher's Linear Discriminant"

Note: this hyperplane separates 2 Gaussians w/
cov $\frac{1}{2} (S_1 + S_2)$

\Rightarrow FLD is optimal when 2 classes are Gaussians
w/ equal covariance matrices.

$$(m_1 - m_2)^2 = [w^T (\mu_1 - \mu_2)]^2 = \underbrace{w^T (\mu_1 - \mu_2) \cdot (\mu_1 - \mu_2)^T w}_{= w^T S_B w}$$