

High Dimensional

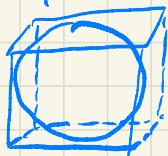
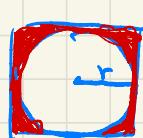
The quality of BDR depends on the CCD estimate

How does it work when X is high-dimensional?

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

Examples

(1) Consider a hypercube and an inside hypersphere in \mathbb{R}^d



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

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

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

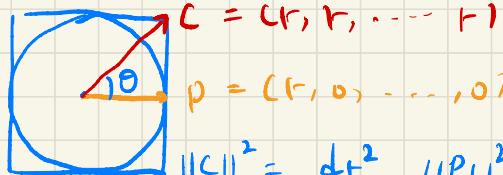
$$\text{Volume of hypercube: } (2r)^d$$

$$\text{Let } f_d = \frac{\text{Volume sphere}}{\text{Volume cube}} = \frac{\pi^{\frac{d}{2}}}{2^d \Gamma(\frac{d}{2} + 1)}$$

$$d = 1 \quad 2. \quad 3 \quad \rightarrow \infty$$

$$f_d = 1 \quad 0.785 \quad 0.524 \quad \rightarrow 0$$

As d increases, the volume of cube's corners increases



$$C = (r, r, \dots, r)$$

$$P = (r, 0, \dots, 0)$$

$$\|C\|^2 = dr^2 \quad \|P\|^2 = r^2$$

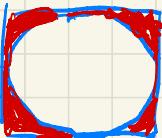
$$\cos \theta = \frac{C^T P}{\|C\| \|P\|} = \frac{r^2}{\sqrt{d} r^2} = \frac{1}{\sqrt{d}}$$

As d increases, then $C \perp P$

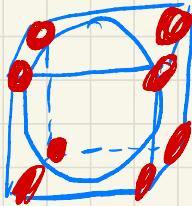
corners are orthogonal to the axis

$$d = 1 \quad \text{---}$$

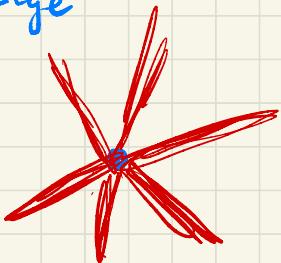
$$d = 2$$



$$d = 3$$

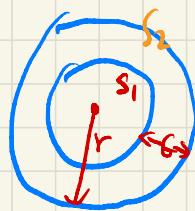


$$d = \text{large}$$



Example 2

Consider a hypersphere Shell of thickness ε



$$V_{\text{shell}} = V(S_2) - V(S_1)$$

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

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

suppose $0 < \varepsilon < r$,

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 rv.

$$\text{Then } E[\|X\|^2] = E(x_1^2 + x_2^2 + \dots + x_d^2)$$

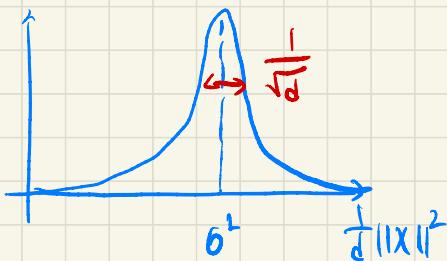
$$= E x_1^2 + E x_2^2 + \dots + E x_d^2$$

$$= E^2 x_1^2 + \text{var}(x_1) + E^2 x_2^2 + \text{var}(x_2) + \dots + E^2 x_d^2 + \text{var}(x_d)$$

$$= d \sigma^2$$

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

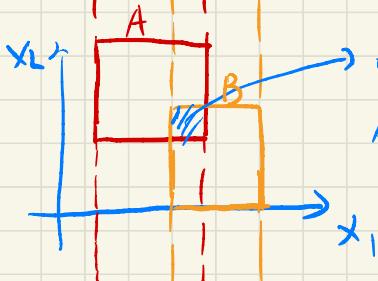
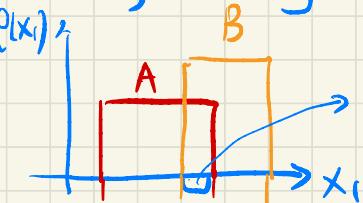
Note: $\|X\|^2$ is a sum of i.i.d. r.v., thus by the central limit theorem, it is concentrated around the mean as $d \rightarrow \infty$, $\frac{1}{d} \|X\|^2 \sim N(b^2, \frac{1}{d})$



In high-dimension, a Gaussian is essentially a shell of radius $b\sqrt{d}$. Most of the density is in the shell (max density is still the mean)

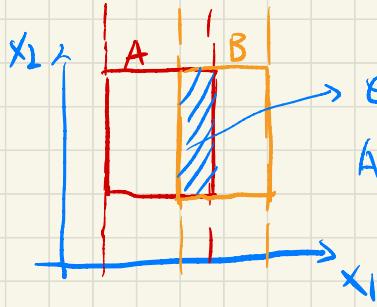
Curse of Dimensionality

In theory, adding new **features** will not increase P_{error}



Adding informative features

$\rightarrow P_{\text{error}} \text{ decreases}$



Errors here!

Adding noninformative features

$\rightarrow P(\text{error})$ is the same
as before.

In practice, for BVR, error increases as feature dimension
increases

Then problem: Quality of the (CD) estimates.

Density estimates in high-dim require more
training samples.

Roughly, desired training size = $O(p^p)$, $p =$
of parameters.

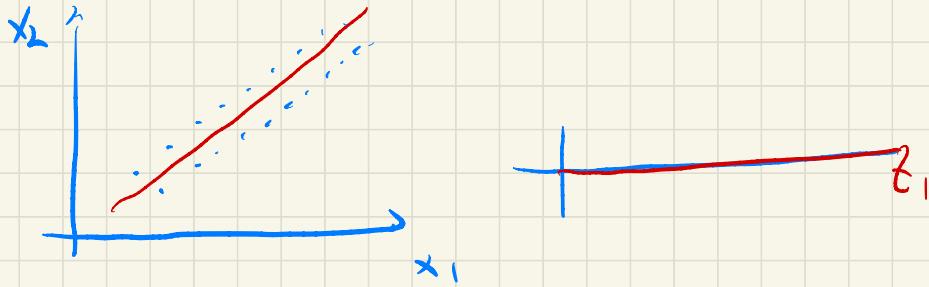
Solution:

- 1) Reduce # of parameters (complexity of model)
(e.g. full covariance \rightarrow diagonal covariance)
- 2) Reduce # of features (dimensionality reduction)
 \rightarrow implicitly reduce # of parameters.
- 3) Create more data
 - a) Bayesian estimate (virtual samples)

b) data augmentation

Linear Dimensionality Reduction

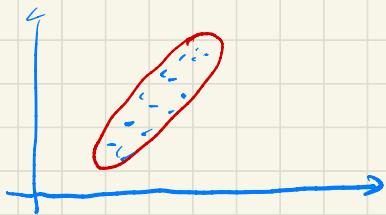
- Summarize correlated features with fewer features
- How to find these correlations?



Correlated data "lives" in a lower-dim subspace with some noise

Principal Component Analysis (PCA)

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



If we fit a Gaussian, it will be "skinny" in some directions.

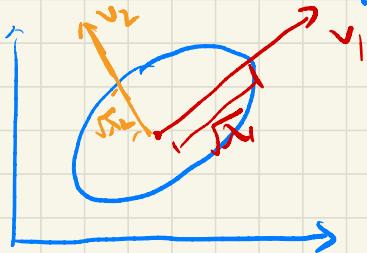
let (v_i, λ_i) be an eigenpair of covariance matrix $\Sigma = V\Lambda V^T$ $V = [v_1, v_2, \dots v_d]$

$$\Lambda = \begin{bmatrix} \lambda_1 & 0 \\ 0 & \ddots & \ddots \\ 0 & & \lambda_d \end{bmatrix}$$

- 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 with larger eigenvalues as "principle components!"



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

1) Calculate Gaussian:

$$\mu = \frac{1}{n} \sum_{i=1}^n x_i, \quad \Sigma = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)(x_i - \mu)^T$$

Training

2) eigen decomposition 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 eigenvalues: $\Phi = [v_1, \dots, v_k]$

5) project new point x onto Φ

$$z = \Phi^T(x - \mu)$$

dim

reduction PCA coefficients

new feature vector, use BDR or other classifiers.

Notes:

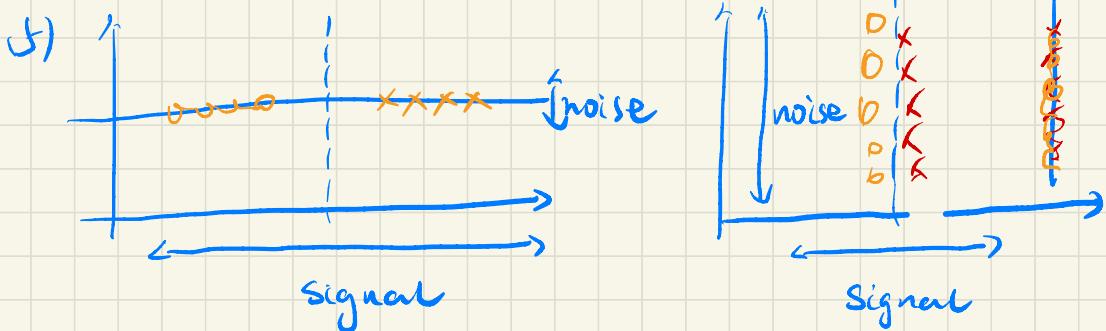
The selection of Φ with $\Phi^T \Phi = I$ also:

- (1) maximize the variance of the projected data
- (2) minimize the reconstruction error of training data
- (3) can be implemented effectively with SVD.

(4) Select k?

- pick k that works in the downstream task (classification)
- pick k to preserve variance of data.

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



Assumption that signal variance is larger than the noise variance.

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

Linear Discriminant Analysis (LDA)

Find the projection that best separates the classes

$$z = w^T x$$

Class statistics :

Class mean

original space

$$\mu_j = \frac{1}{n} \sum_{x_i \in C_j} x_i$$

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 between projected means.

$$(m_1 - m_2)^2 = (w^T (\mu_1 - \mu_2))^2$$

Problem : w is unconstructed \rightarrow need normalization.

Fisher's Idea:

between-class scatter



$$w^* = \underset{w}{\operatorname{argmax}} \frac{(m_1 - m_2)^2}{S_W + S_B} = \underset{w}{\operatorname{argmax}} \frac{w^T S_B w}{w^T S_W w}$$

with-class scatter.

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

$$S_B = (m_1 - m_2)(m_1 - m_2)^T$$

$$S_W = S_1 + S_2$$

Note this hyperplane separates 2 Gaussian with
 $\text{cov} \frac{1}{n}(S_1 + S_2)$

\Rightarrow FLD is optimal when 2 classes are
Gaussian with equal covariance matrices.