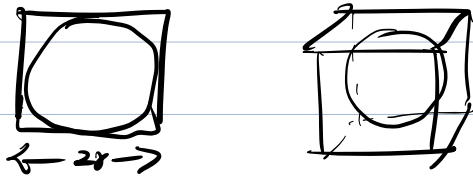


The quality of BDR depends on the CCD estimates
How does it work when \mathcal{X} is high-dimensional?

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

Examples:

c1) Consider a hypercube & an inscribed hypersphere in \mathbb{R}^d



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

Gamma function

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

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

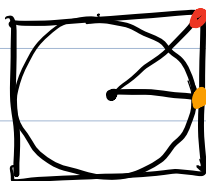
Volume of hypercube $(2r)^d$

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

d	1	2	3	∞
f_d	1	0.785	0.524	0

factorial overweights

As d increases, the volume of the corner increases.



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

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

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

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

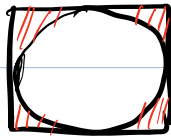
$$\cos\theta = \frac{c^T p}{\|c\| \|p\|} = \frac{r^2}{\sqrt{d} r^2} = \frac{1}{\sqrt{d}}$$

as $d \nearrow \Rightarrow \cos\theta = 0 \Rightarrow c \perp p!$
(corner is orthogonal to axis)

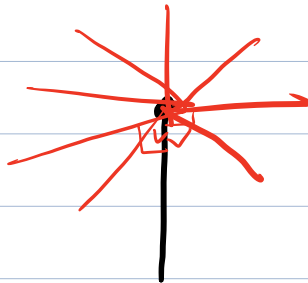
$d=1$



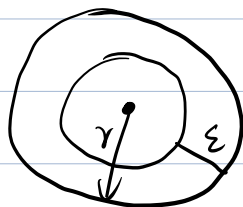
$d=2$



$d \Rightarrow \infty$



Example 2 a hypersphere shell of thickness ε



$$U_{\text{shell}} = U(s_2) - U(s_1)$$

$$= \left(1 - \frac{U(s_1)}{U(s_2)}\right) \cdot U(s_2)$$

$$\frac{U(s_1)}{U(s_2)} = \frac{C(r-\varepsilon)^d \pi^{d/2}}{\frac{r^d \pi^{d/2}}{\Gamma(\frac{d}{2}+1)}} = \left(1 - \frac{\varepsilon}{r}\right)^d$$

Suppose $0 < \varepsilon < r$

$$\text{as } d \nearrow, \frac{U(s_1)}{U(s_2)} \rightarrow 0$$

$$U_{\text{shell}} \rightarrow U(s_2)$$

All the volume is in the shell of the hypersphere

Example 3 high-dim Gaussian

$$\text{let } x \sim \mathcal{N}(0, \sigma^2 I)$$

i.e. $x_i \sim \mathcal{N}(0, \sigma^2)$ i.i.d. r.v.

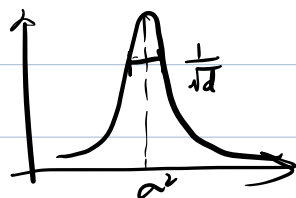
$$\text{Then } E[\|x\|^2] = E[x_1^2 + x_2^2 + \dots + x_d^2] = d\sigma^2$$
$$E[\frac{1}{d}\|x\|^2] = \sigma^2$$

Note $\|x\|^2$ is a sum of i.i.d. r.v.

By central limit theorem, it is concentrated

around mean as $d \rightarrow \infty$

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



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

why $\|\frac{1}{\sqrt{d}}x\|^2$

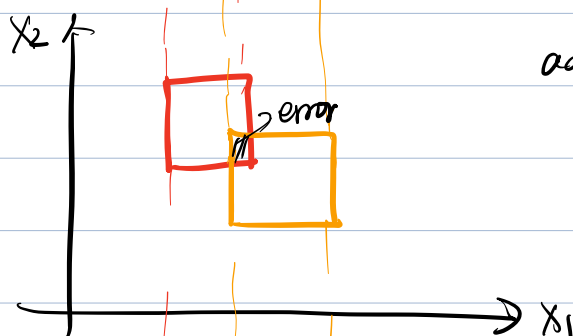
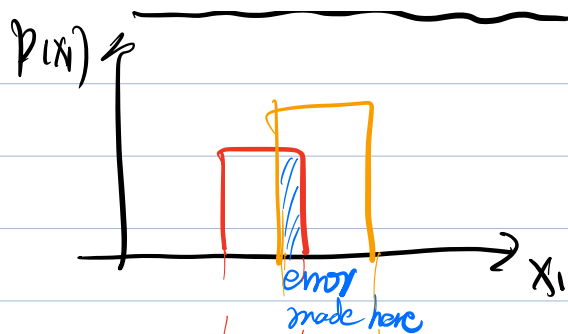
However, the mean is still 0

though density is in the shell

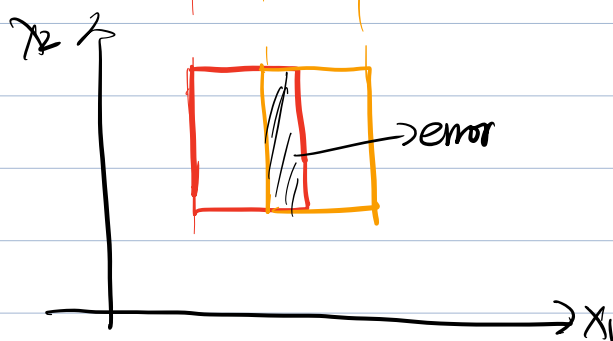
Our intuition in 2d
 \Rightarrow density should
be around the mean

Curse of dimensionality

In theory, add a feature will not
increase the error ✖



adding informative features
 $\rightarrow p(\text{error})$ decreases



Adding non-informative
 features $p(\text{error})$ is
 the same as
 before

In practice, for BQR, error increases as the
 feature dim increases

The problem is the quality of the CCD estimate.

Density estimates in high-dim require more
 training samples.

Roughly, desired training set size $= O(e^P)$

$P = \#$ of parameters

Solution:

→ things we optimize

1) Reduce # of parameters (complexity of model)

eg. full cov \Rightarrow diag cov

2) Reduce # of features (dimensionality reduction)

→ implicitly reduce # of parameters

3) Create more data

(a) Bayesian estimation (virtual samples)

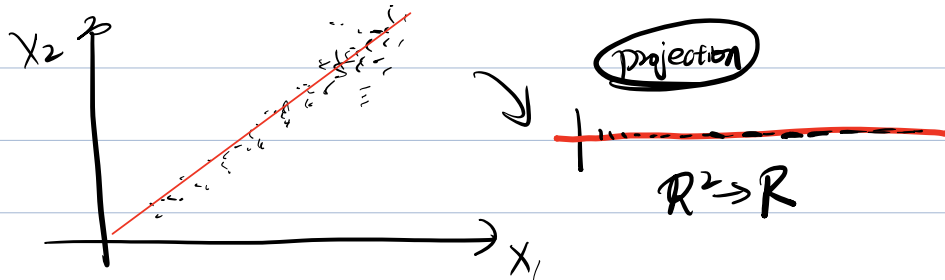
(b) data augmentation.

(eg. $\boxed{A} \Rightarrow \boxed{A}$ add noise
 $\boxed{A} \Rightarrow \boxed{A}^{\circ}$)

Linear Dimensionality Reduction

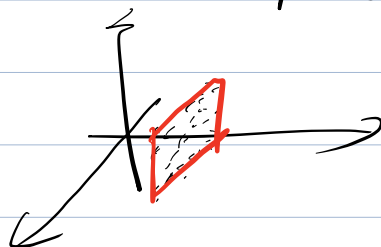
- Summarize correlated features w/ fewer features.

- How do we find these correlations?



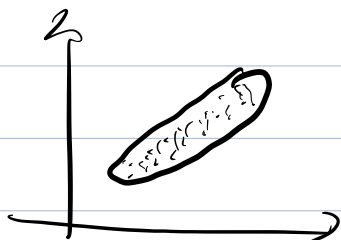
Correlated data "lives" in a lower-dim

subspace (w/ some noise)



PCA (Principle Component Analysis)

Idea If the data can fit into a subspace,
then it should be flat in full space.



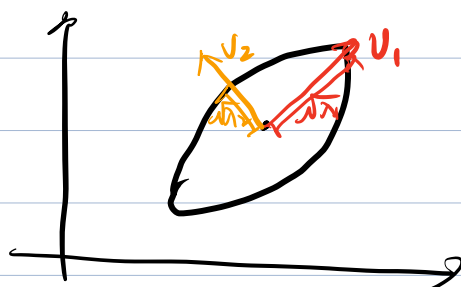
If we fit a Gaussian,
it will be "skinny" in
some directions.

Let (v_i, λ_i) be an eigenpair of covariance
matrix Σ

$$\Sigma = U \Lambda U^T, \quad U = [v_1, v_2, \dots, v_d]$$

$$\Lambda = \begin{bmatrix} \lambda_1 & & 0 \\ & \lambda_2 & \\ 0 & & \ddots \end{bmatrix}$$

- 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/ large eigenvalues
as "principle components".

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

1) Calculate Gaussian.

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

- training
- 2) eigen decomposition $\Sigma = V \Lambda V^T$
 - 3) order the eigenvalues $\lambda_1 > \lambda_2 > \dots > \lambda_d > 0$
 - 4) select top k eigenvectors. $\Phi = [v_1, \dots, v_k]$
- dim reduction.
- 5) project new points x onto Φ
- $$z = \Phi^T (x - \underbrace{\bar{x}}_{\text{mean}}) \leftarrow \text{PCA coefficients. (row feature vectors)}$$
- Central the data

Notes:

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

- ps
- (1) maximize the variance of the projected data into z
 - (2) minimize the reconstruction error of training data



- (3) Can be implemented efficiently w/ SVD ps 7-8
- (4) How to select k ?