

Dimensionality

The quality of BDR depends on the class conditional densities (CCD) estimates.

• How does it work in high-dimension?

"High dimensional spaces are weird."

Do not trust your intuition.

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



$d=2$



$d=3$

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

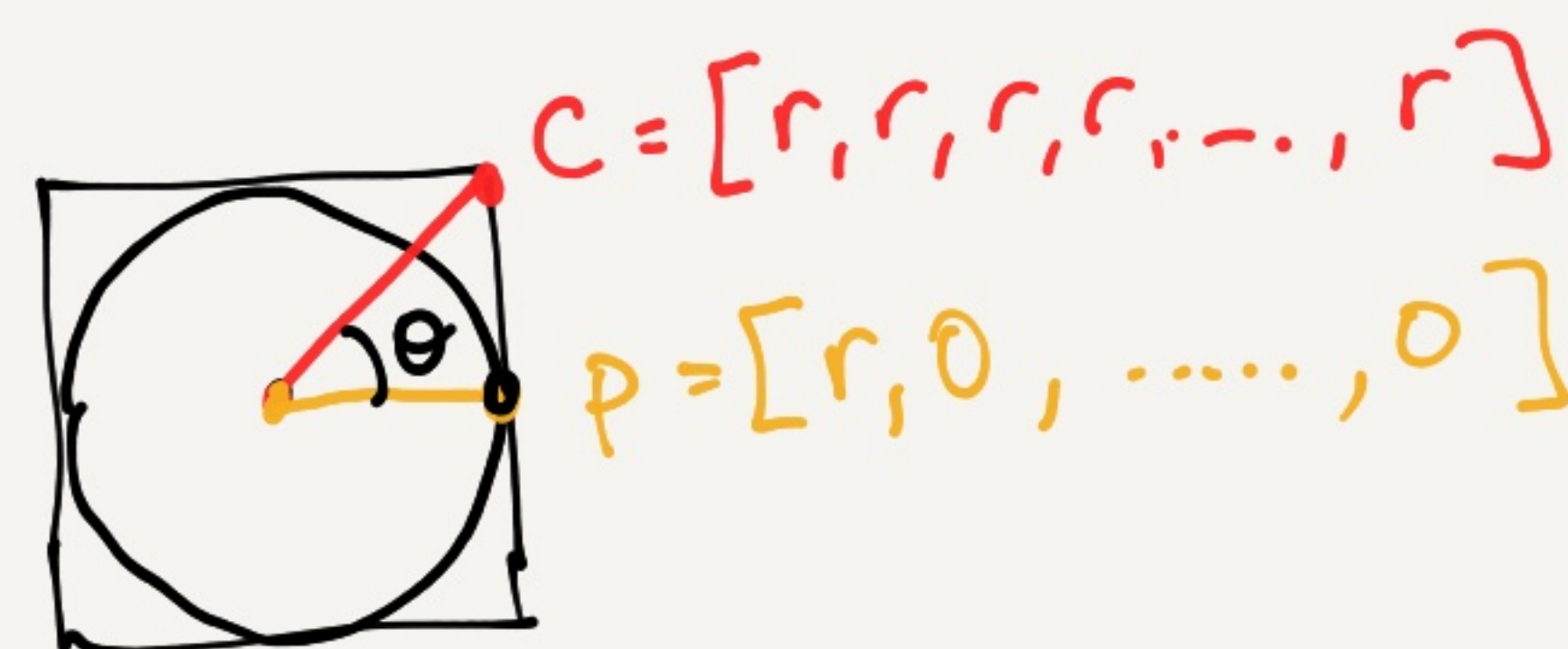
$\Gamma(n) = \int_0^\infty e^{-x} x^{n-1} dx$ (Gamma function)
 \approx factorial for real numbers.
 $\Gamma(n+1) = n!$

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

Let $f_d = \left[\frac{\text{volume of h-sphere}}{\text{volume of h-cube}} \right] = \frac{\pi^{d/2}}{2^d \Gamma(\frac{d}{2} + 1)}$

as d increases, $f_d \rightarrow 0$.

\therefore the volume of the corners increases.



$$\|c\|^2 = d r^2$$

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

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

Assignment Project Exam Help

<https://powcoder.com>

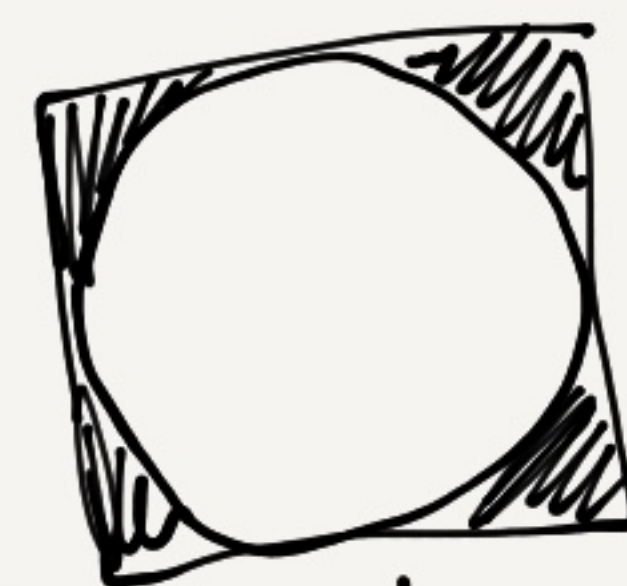
Add WeChat powcoder

As d increases, then $\cos \theta \rightarrow 0$
 i.e. $c \perp p$ (corner is orthogonal to the axis)

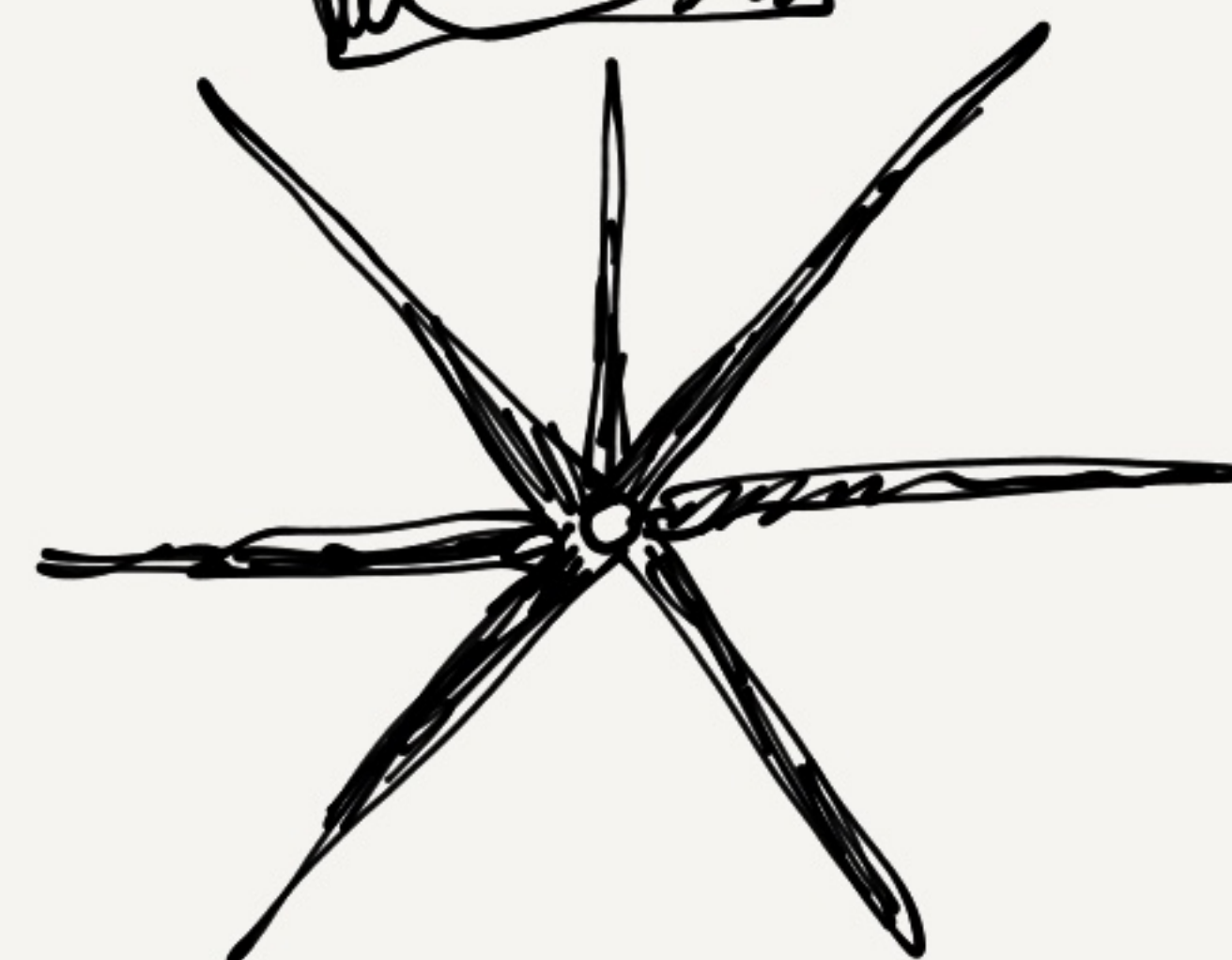
$d=1$



$d=2$

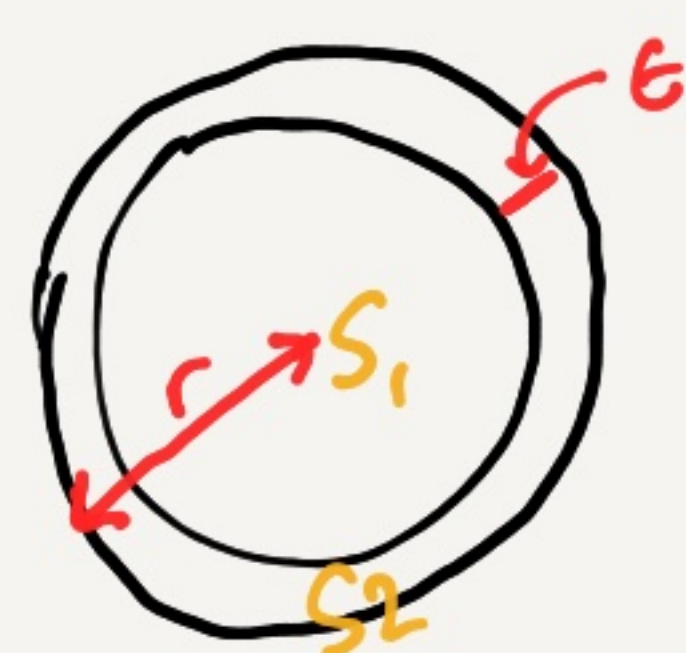


large d



all the volume in the corner "spikes"

Example 2) consider a hypersphere shell of thickness ϵ .



Volume of hypershell

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

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

$$\frac{V(S_1)}{V(S_2)} = \dots = \left(1 - \frac{\epsilon}{r}\right)^d$$

< 1

Suppose $0 < \epsilon < r$: As d increases, $\frac{V(S_1)}{V(S_2)} \rightarrow 0$

Thus: $V_{\text{shell}} \rightarrow V(S_2)$

"all the volume is in the shell!!!" (?)

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Example 3: high-dim Gaussian.

m.v. Gaussian $X \sim N(0, \sigma^2 I_d)$

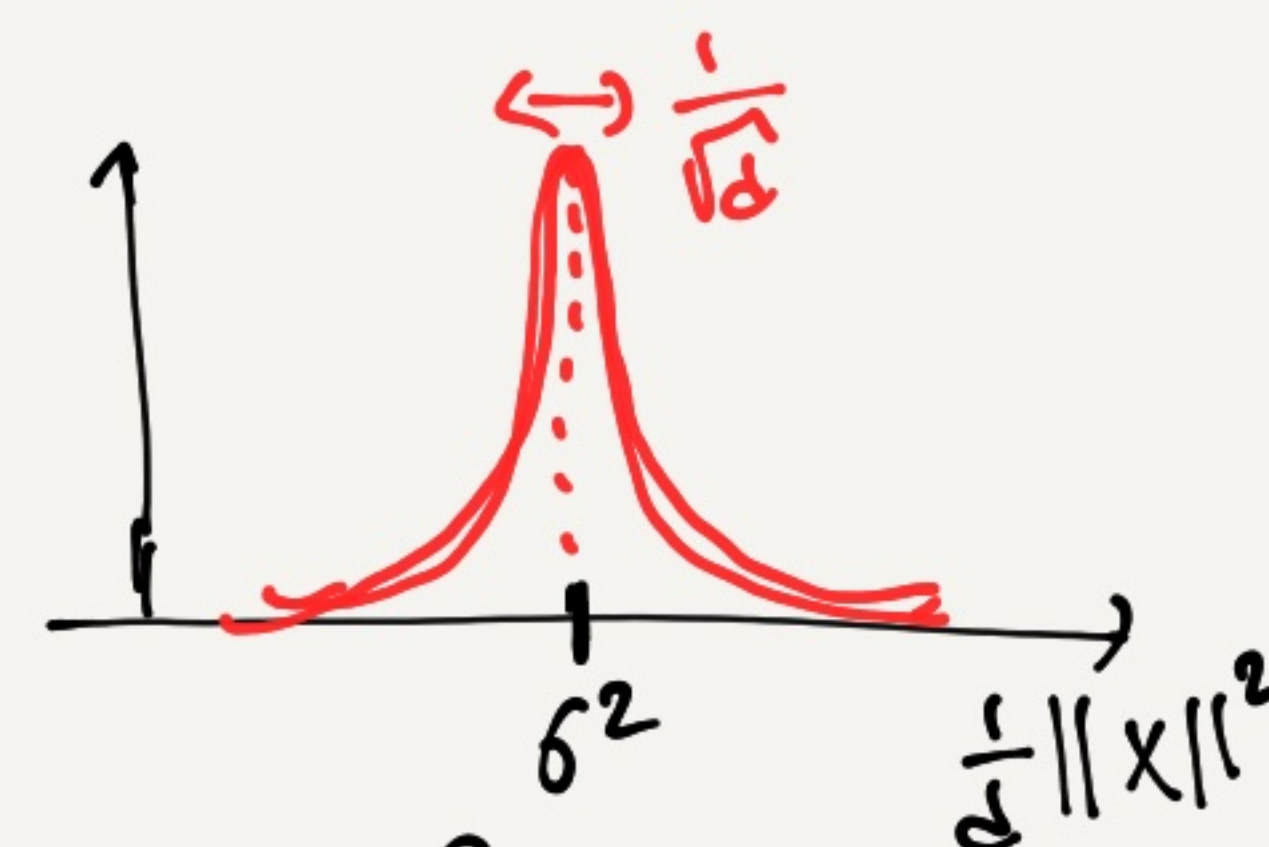
thus $x_i \sim N(0, \sigma^2)$ (iid r.v.)

Then: $E\left[\frac{1}{d}\|X\|^2\right] = \frac{1}{d} E\left[\underbrace{x_1^2 + x_2^2 + \dots + x_d^2}_{\substack{\text{d terms} \\ \sigma^2}}\right] = \sigma^2$

Note: $\|X\|^2 = \sum_i x_i^2 \rightarrow$ sum of iid r.v.
 $\underbrace{\quad}_{\substack{\text{r.v.} \\ \text{Sum of r.v.}}}$

So by the central limit theorem, $\|X\|^2$ is concentrated around the mean as $d \rightarrow \infty$

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



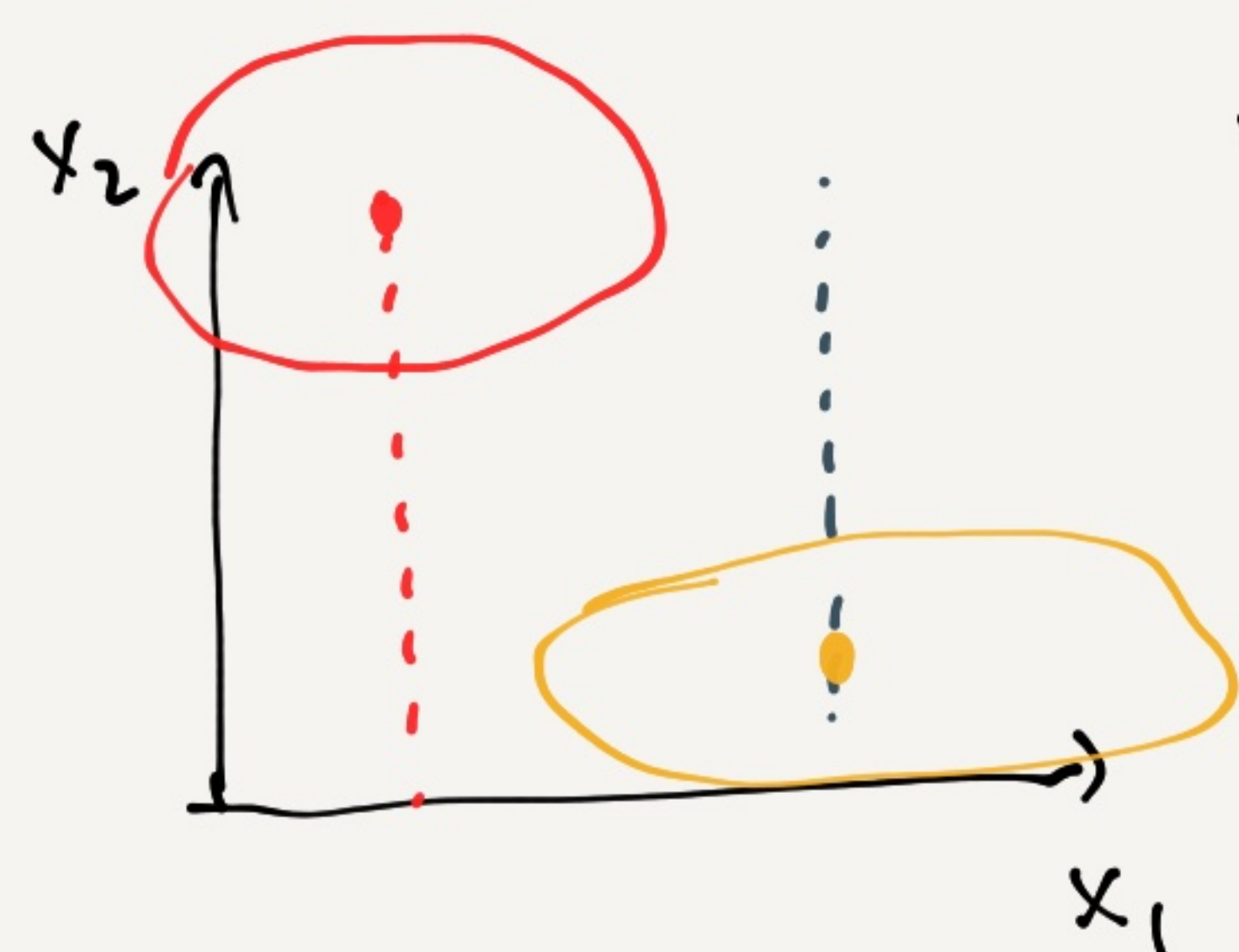
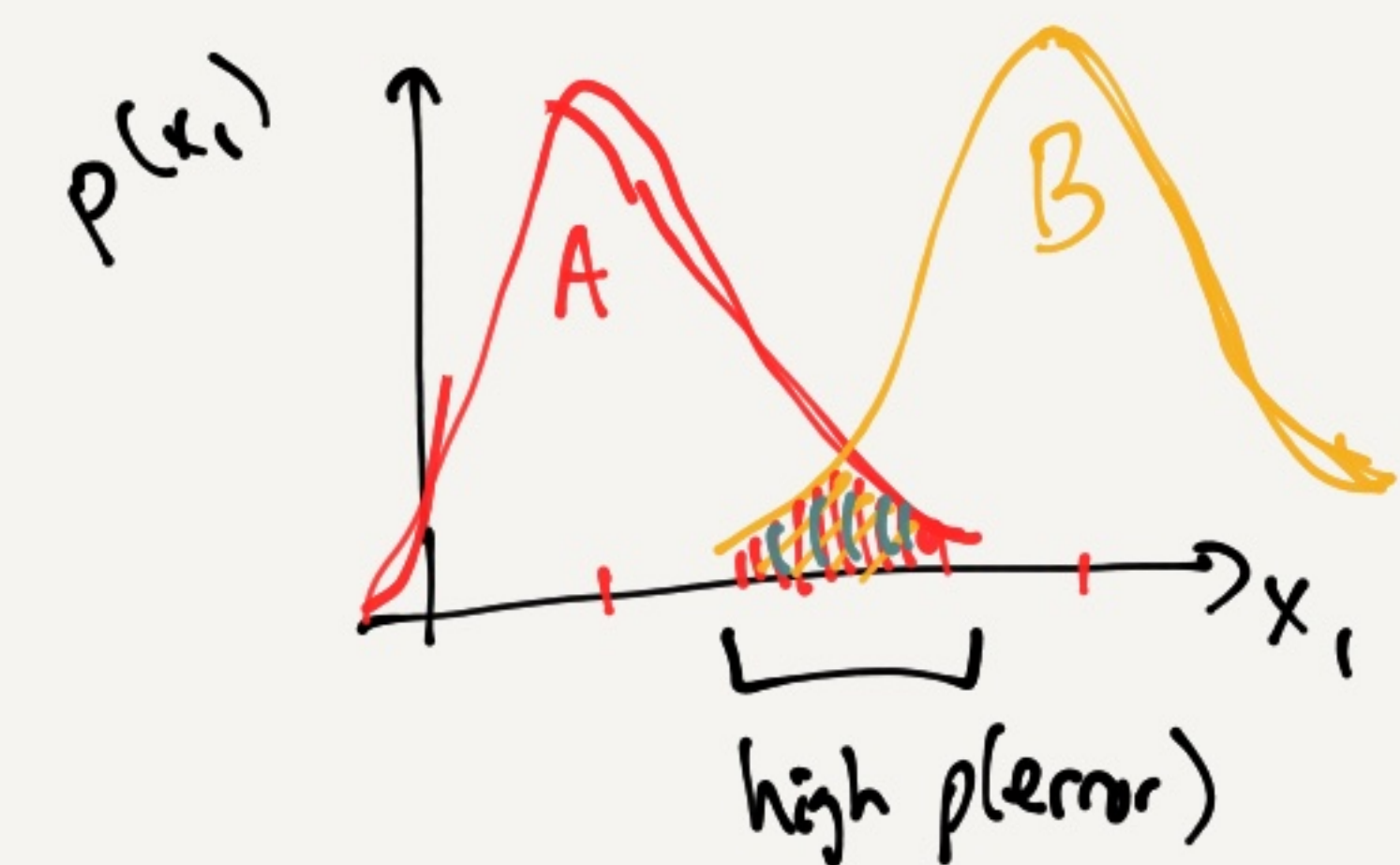
In high-dim, a Gaussian

is a shell of radius $\sigma\sqrt{d}$ — most of the distribution is inside the shell.

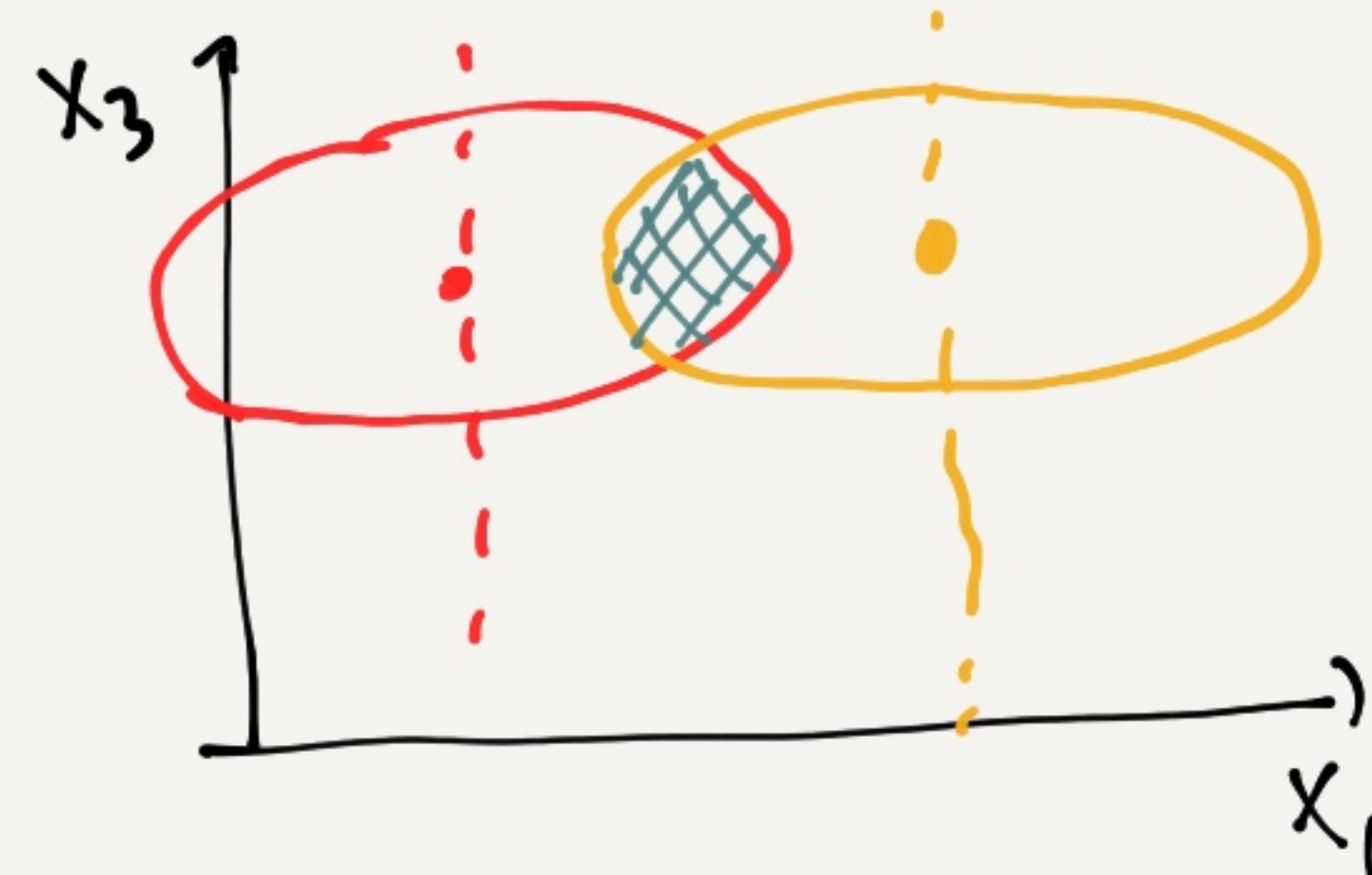
(the max density value is still the mean)

Curse of dimensionality

In theory, adding new features will not hurt the classifier (i.e. increase $p(\text{error})$)



x_2 is an informative feature
 $\rightarrow p(\text{error})$ decrease



x_3 is uninformative,
 $\rightarrow p(\text{error})$ is the same as before.

But: in practice, for BDR the error increases as we add more features (add dimensions to the input)

The problem: quality of the (CD) estimates.
density estimates in high-dim require more training data!

e.g. histogram in d-dim over unit cube $[0,1]^d$
 $\Rightarrow 10$ bins per dimension $\Rightarrow 10^d$ bins overall
to just have 1 sample per bin,
we need 10^d samples \rightarrow increases exponentially w/ the dimension.

Solutions: 1) Reduce # parameters (complexity of model)
e.g. full cov \rightarrow diagonal cov.

✓ 2) Reduce # of features (dimensionality reduction)
 \Rightarrow implicitly reduce # of parameters.

3) "Create" more data

a) virtual samples (Bayesian estimation)

b) data augmentation (apply xforms to the data that keep the class the same)

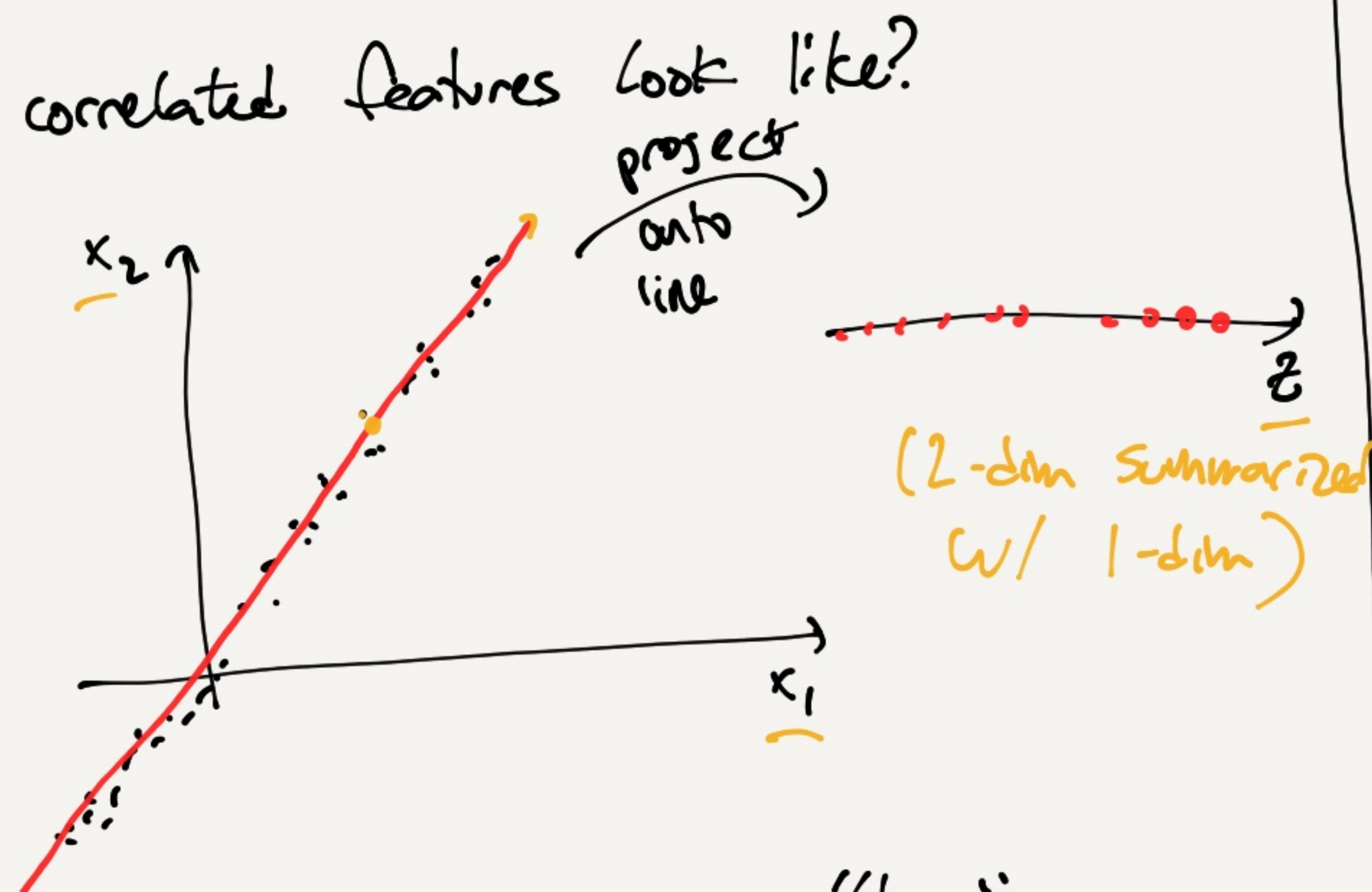
Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Linear Dimensionality Reduction

- summarize correlated features w/ a set of fewer features.
- What do correlated features look like?



- the correlated features live on a "line" or subspace (w/ some noise)

Assignment Project Exam Help

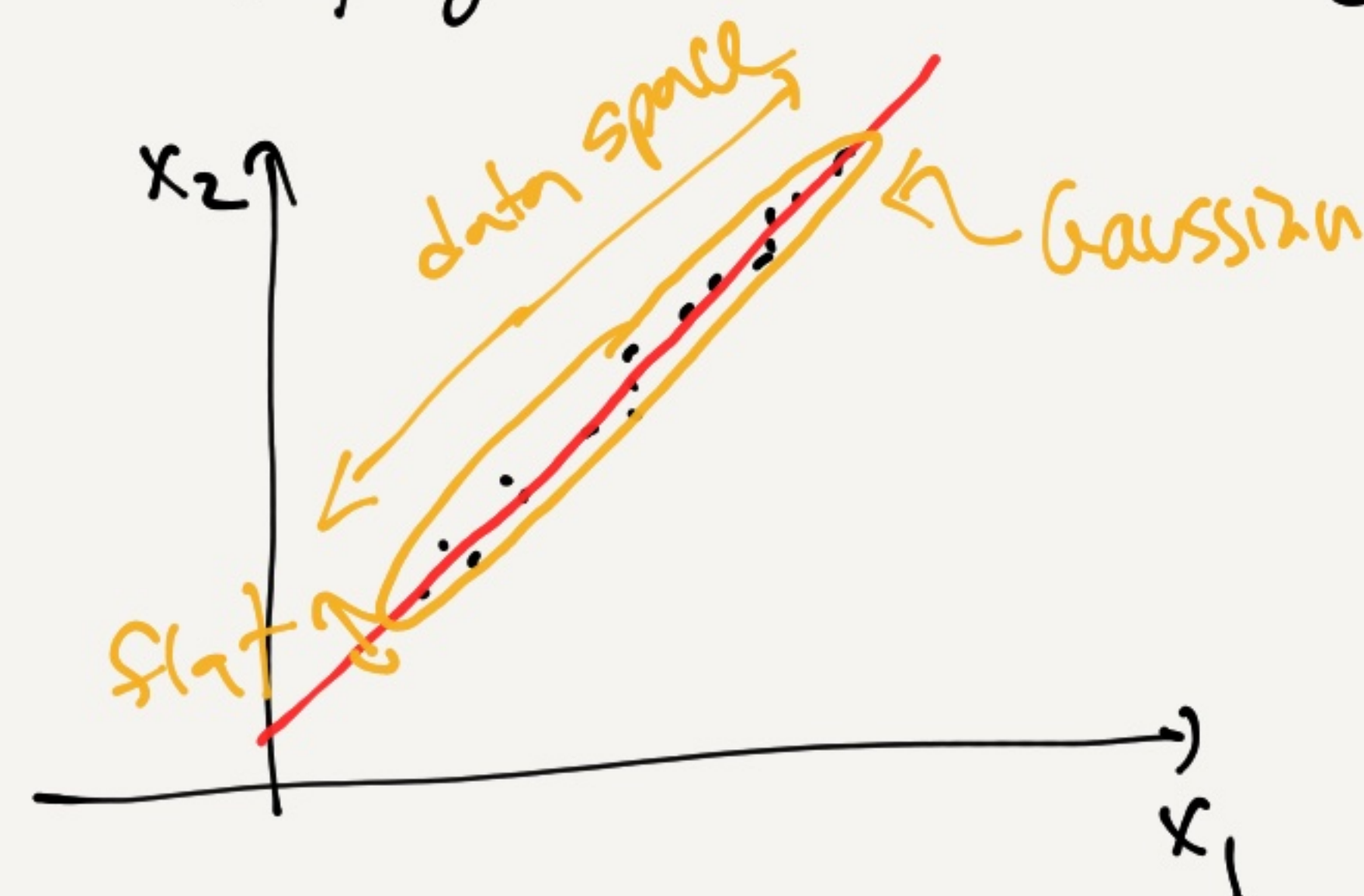
<https://powcoder.com>

Add WeChat powcoder

Principal Component Analysis (PCA)

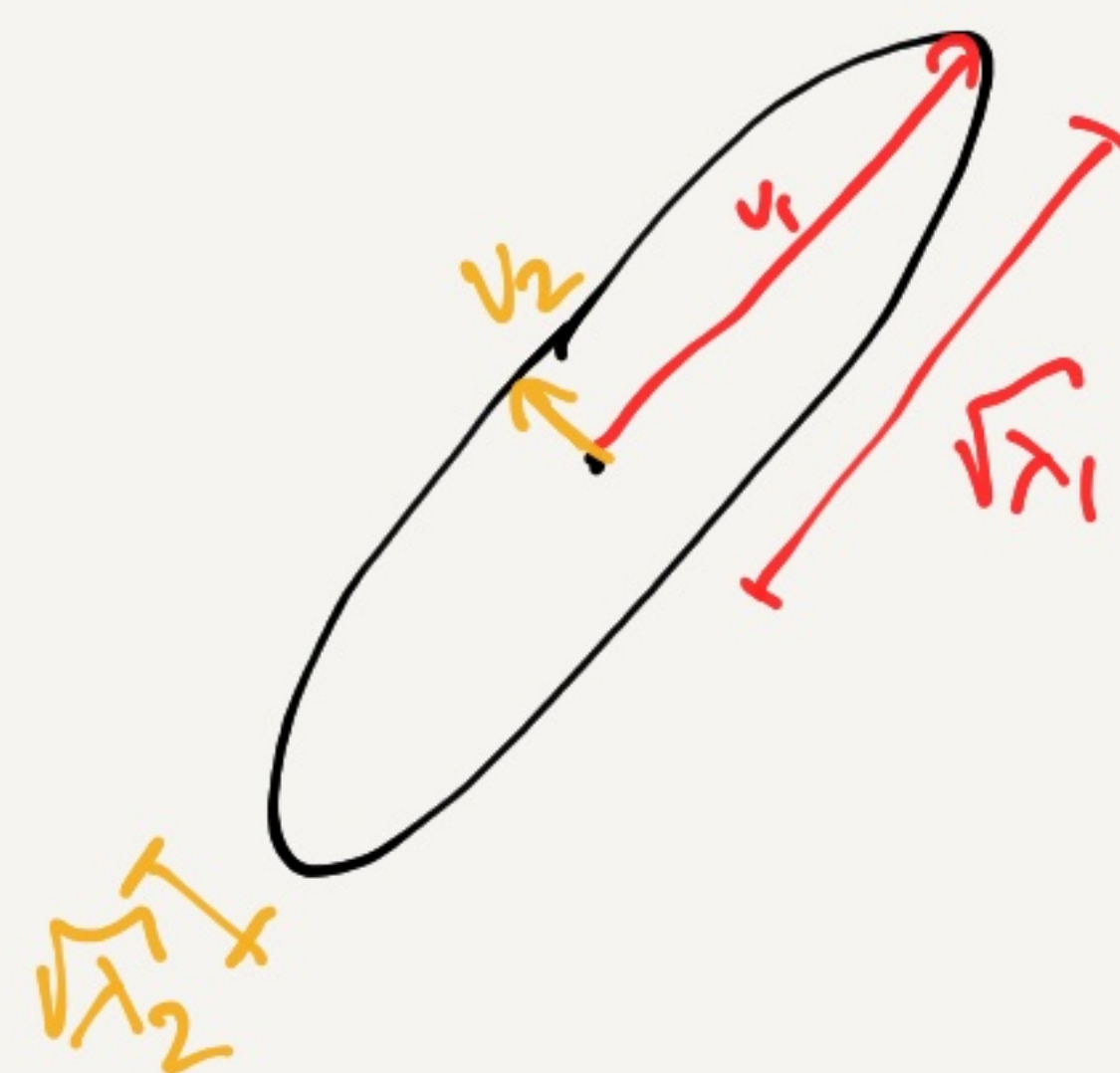
Idea: if the data lives in a low-dim subspace, then it should be "flat" in the full space.

\Rightarrow if we fit a Gaussian \rightarrow it will be highly skewed (skinny ellipses)



- let (v_i, λ_i) be an eigenvector / eigenvalue pair of $\text{cov } Z$.

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



Thus the eigenvalues tell us which directions are flat.

\Rightarrow select v_i w/ largest eigenvalue for the projection

PCA: Given dataset $\{x_1, \dots, x_N\}$ of dim k .

learning

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 decomposition of $\Sigma = V \Lambda V^T$

$V = [v_1 \dots v_d]$ = matrix of eigenvectors

$\Lambda = \begin{bmatrix} \lambda_1 & 0 \\ & \ddots \\ 0 & \lambda_d \end{bmatrix}$ = diagonal matrix of eigenvalues.

3) Sort eigenvalues: $\lambda_1 > \lambda_2 > \lambda_3 > \dots \geq 0$

4) select top- k eigenvectors: $\Phi = [v_1 \dots v_k]$

Principal components

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

5) project point x onto Φ space:

$$z = \Phi^T (x - \mu) \leftarrow \text{"PCA coefficients"}$$

z is the new feature, apply BDR as before

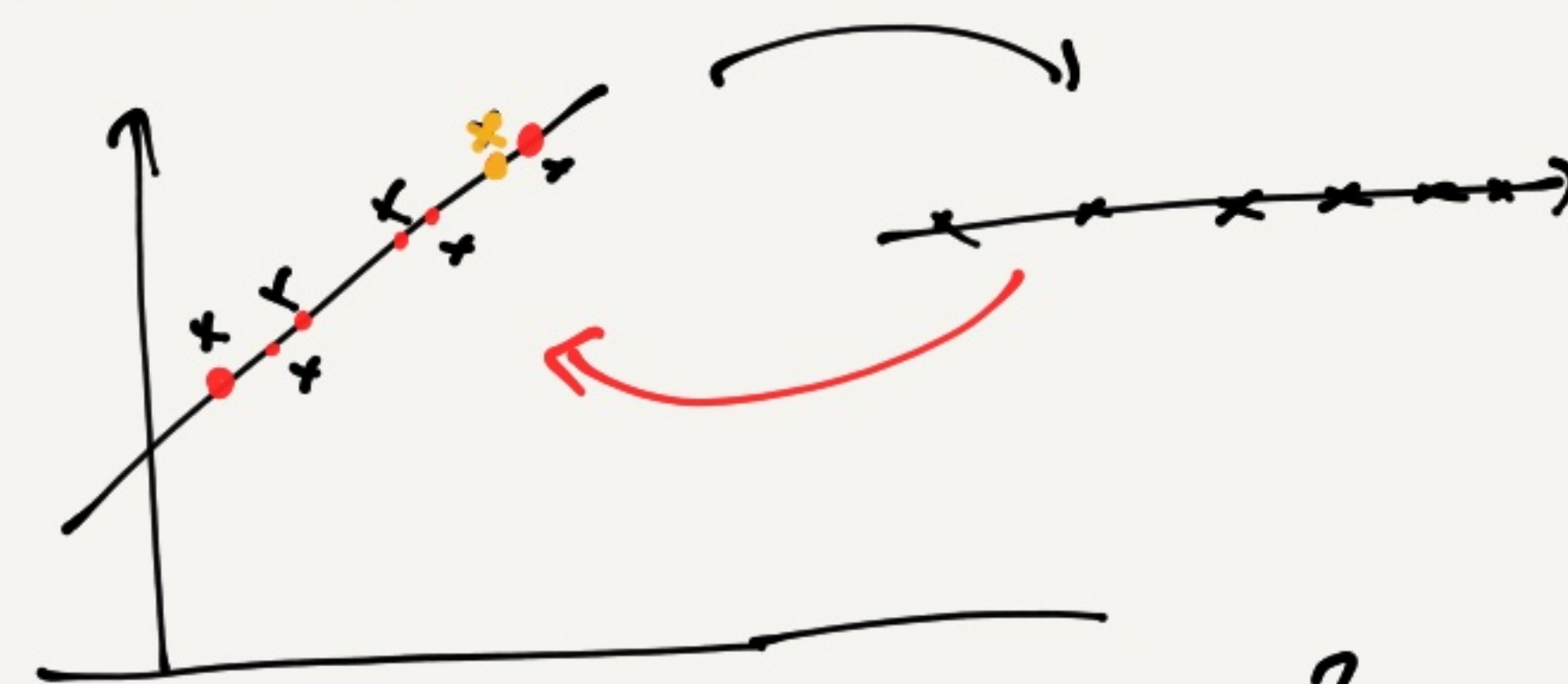
dim. reduction

Notes:

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

1) maximize the variance of the projected training data
i.e. $\sum_i \|z_i\|^2$ (PS 7.3)

2) minimizes the reconstruction error of the training data
PS 7.2



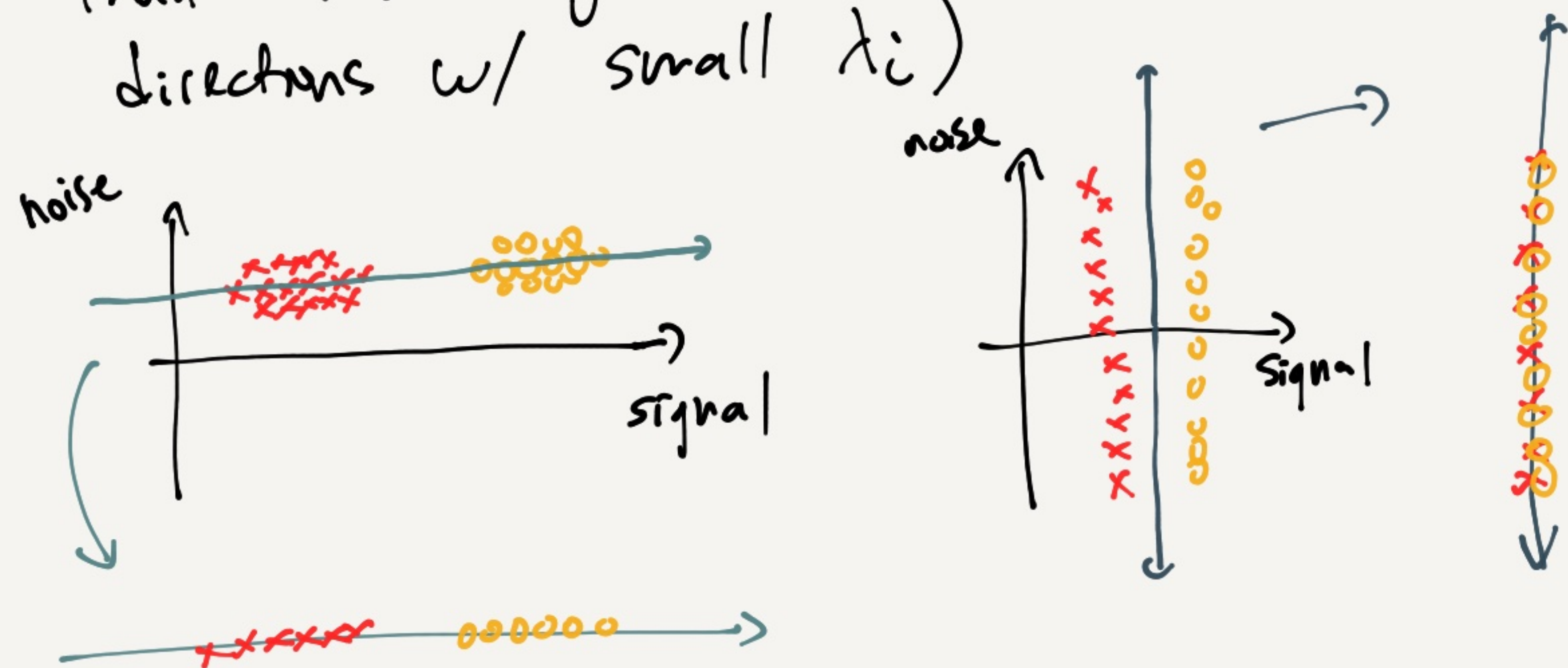
3) for high-dim data, can use SVD for an efficient implementation. (PS 7-4)

4) Selecting k : 1) pick a value of k (that works for your problem)

2) pick k to preserve $p\%$ of the variance

$$p = \frac{\sum_{i=1}^k \lambda_i}{\sum_{i=1}^d \lambda_i} \left[\begin{array}{l} \leftarrow \text{variance preserved} \\ \leftarrow \text{total variance} \end{array} \right]$$

5) assume the "noise" variance is smaller than the "signal" variance. (throw away directions w/ small λ_i)



\Rightarrow has consequences on classification.

6) PCA is optimal for representation (preserving variance)

but not necessarily classification.

- there is no way to fix this because we don't have the classes!

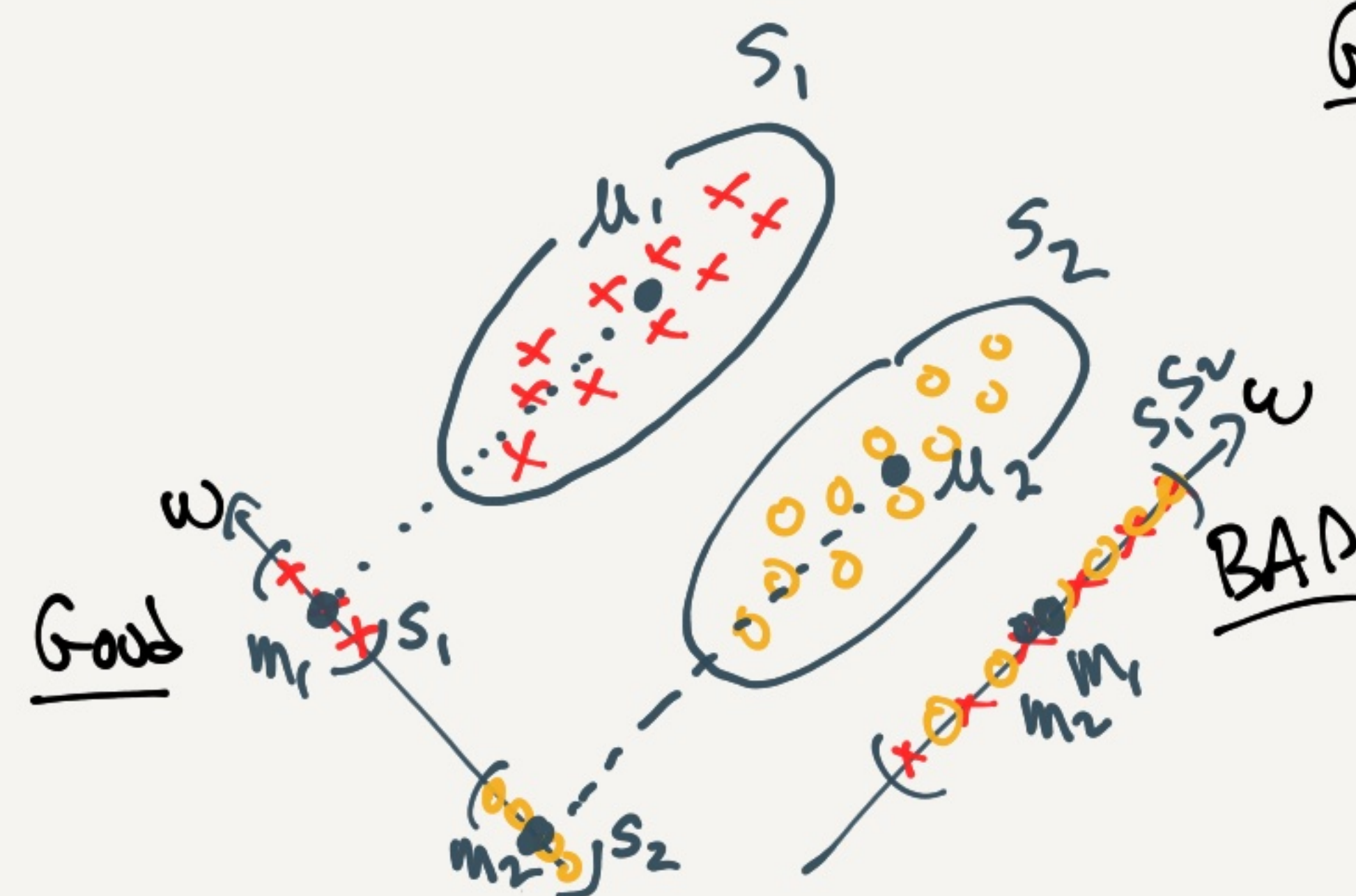
\Rightarrow Fix this using the class information when doing dim. reduction.

Fisher's Linear Discriminant (FLD)

(Linear discriminant analysis (LDA))

Goal: find the projection that best separates the classes.

$$z = w^T x$$



Class Statistics

original space

class mean: $\mu_j = \frac{1}{N_j} \sum_{x_i \in C_j} x_i$

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

1-d space

$$m_j = w^T \mu_j$$

$$S_j = w^T S_j w$$

Idea: maximize the distance between the projected means: $(m_1 - m_2)^2 = (w^T(\mu_1 - \mu_2))^2$
problem: w is unconstrained \rightarrow need normalization.

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Fisher's Idea

$$w^* = \operatorname{argmax}_w$$

$$\frac{\overbrace{(\mu_1 - \mu_2)^2}^{\text{between-class scatter}}}{\underbrace{S_1 + S_2}_{\text{within-class scatter}}}$$

$$= \operatorname{argmax}_w \frac{w^T S_B w}{w^T S_W w}$$

$$\begin{cases} S_B = (\mu_1 - \mu_2)(\mu_1 - \mu_2)^T \\ S_W = S_1 + S_2 \end{cases}$$

⋮
↓
generalized eigenvalue problem (tutorial)

← Fisher's Linear Discriminant

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

Note: this also defines the decision boundary of a Gaussian classifier w/ cov $\frac{1}{N}(S_1 + S_2)$.

⇒ FLD is optimal when the 2 classes are Gaussian w/ equal covariance.

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder