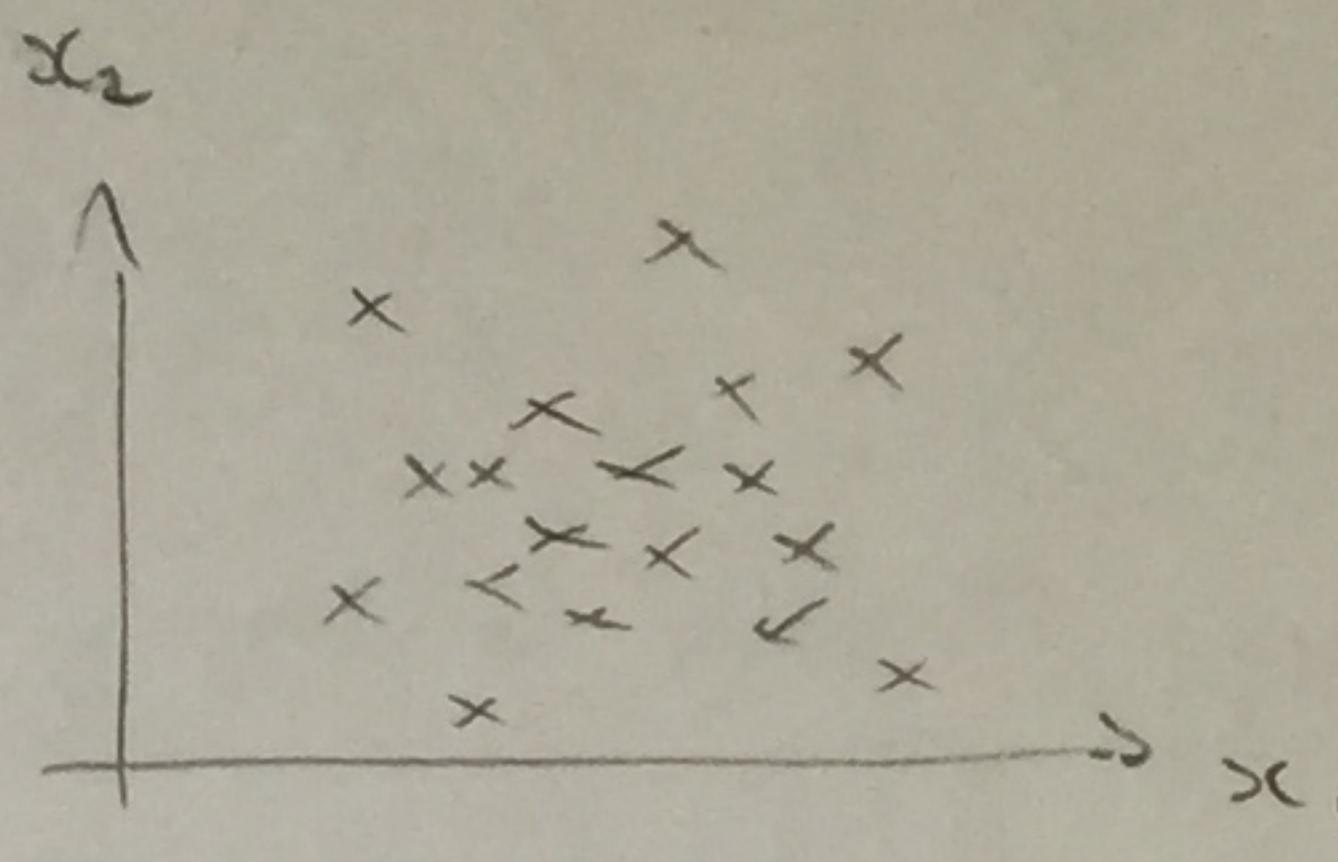


Anomaly Detection:

e.g. Aircraft mfg: $x_1 = \text{heat gen.}$
 $x_2 = \text{vib. intensity} \Rightarrow$



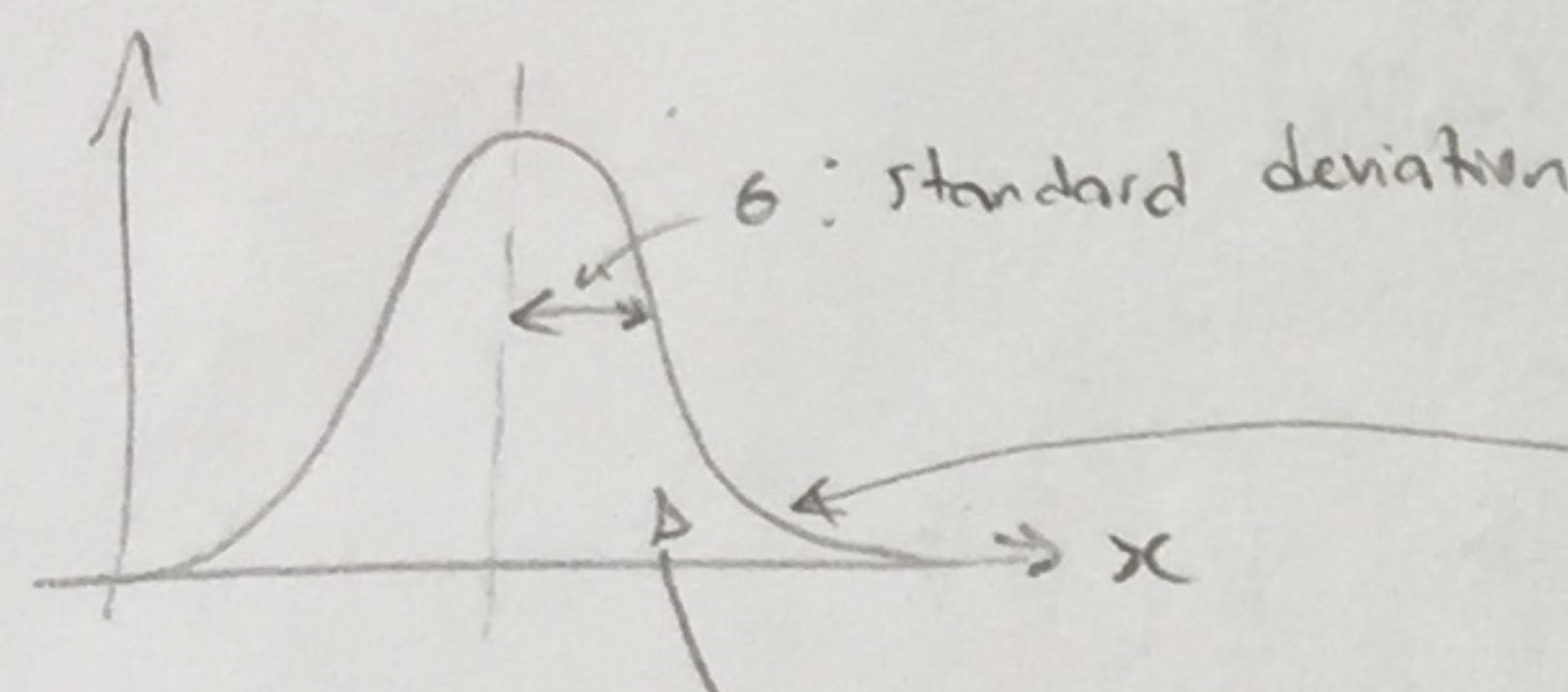
↳ how to determine what constitutes an anomaly above?

\Rightarrow can build model $p(x) \rightarrow$ if $p(x_{\text{test}}) < \epsilon \rightarrow$ flag anomaly!

Applications: fraud detection, manufacturing

Gaussian Distribution Review:

$x \in \mathbb{R}$, distributed gaussian w/ mean μ & variance $\sigma^2 \rightarrow x \sim N(\mu, \sigma^2)$ distributed as normal



$$p(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi} \sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

b/c probability dist., area = 1

\Rightarrow given $x^{(i)} \sim N(\mu, \sigma^2)$, can guess μ & σ^2 by a sample dataset:

$$\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)} \quad \sigma^2 = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu)^2 \leftarrow \begin{array}{l} \text{max likelihood estimations} \\ \text{sometimes } (m-1) \end{array}$$

Anomaly Detection Algorithm

$$p(x_i; \mu_i, \sigma_i^2)$$

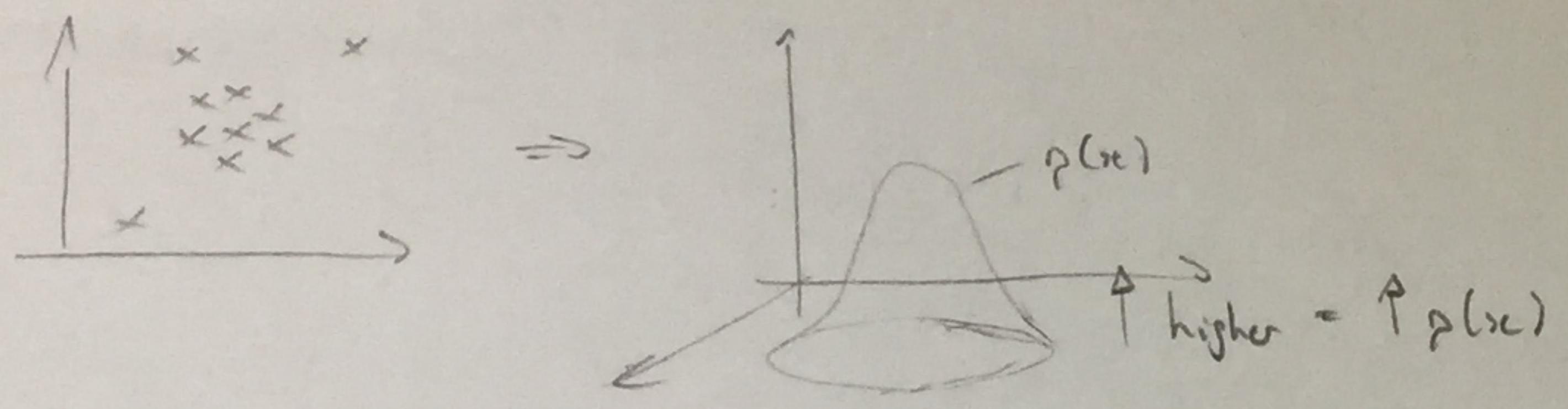
$$p(x) = p(x_1) \cdot p(x_2) \cdot p(x_3) \cdots p(x_n)$$

$$= \prod_{j=1}^n p(x_j; \mu_j, \sigma_j^2)$$

① Choose features x_i indicative of anomalous examples

② Find μ_j , σ_j^2 for $j=1-n$ features

$$\textcircled{3} \quad p(x) = \prod_{j=1}^n p(x_j; \mu_j, \sigma_j^2) = \prod_{j=1}^n \frac{1}{\sqrt{2\pi} \sigma_j} \exp\left(-\frac{(x_j - \mu_j)^2}{2\sigma_j^2}\right) \rightarrow \text{if } p(x) < \epsilon, \text{ anomaly!}$$



⇒ making decisions in evaluating a learning algorithm always better if we have some metric

↳ eg we have labeled data (anomaly, non-anomaly)

↳ extract • training set (assume all non-anomalies)

- cross val set
- test set

eg 10 000 normal, 20 anomalous →

train:	6000 normal
cv:	2000 normal, 10 anom.
test:	2000 normal, 10 anom

① Fit model $p(x)$ on train

② Predict $y=1$ (anomaly), $y=0$ (normal) on C.V.

③ Use evaluation metric → • True Pos, etc.

• Precision / Recall

④ Evaluate on test
• F₁-Score

→ Use CV to choose ϵ

⇒ Why not just use a supervised learning algorithm?

↳ usually only have small number of $y=1$ (anomaly) examples relative to $y=0$!

↳ future anomalies may look nothing like current anomalies

↳ use supervised if we have lots of anomaly examples (eg spam)

ANOMALY:

- fraud
- mfg
- machines in data centres

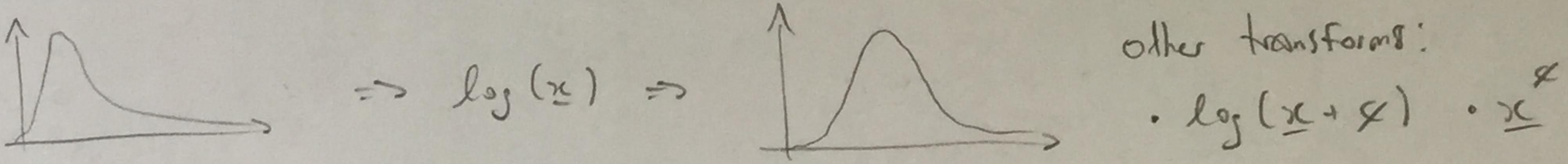
SUPERVISED LEARNING:

- spam
- weather
- cancer prediction

Choosing the Right Features for Anomaly Detection:

⇒ plot histogram of data to make sure \underline{x} is gaussian

↳ transform data to look gaussian (e.g. $\log(\underline{x})$) if not gaussian!

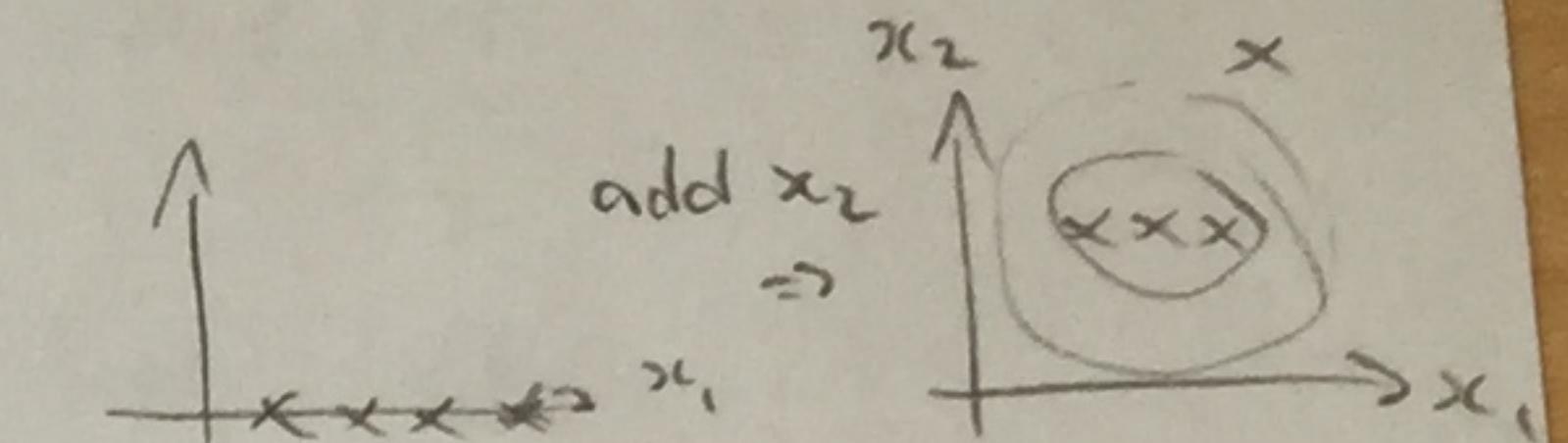


⇒ Want $p(x)$ large for normals, $p(x)$ small for anomalies!

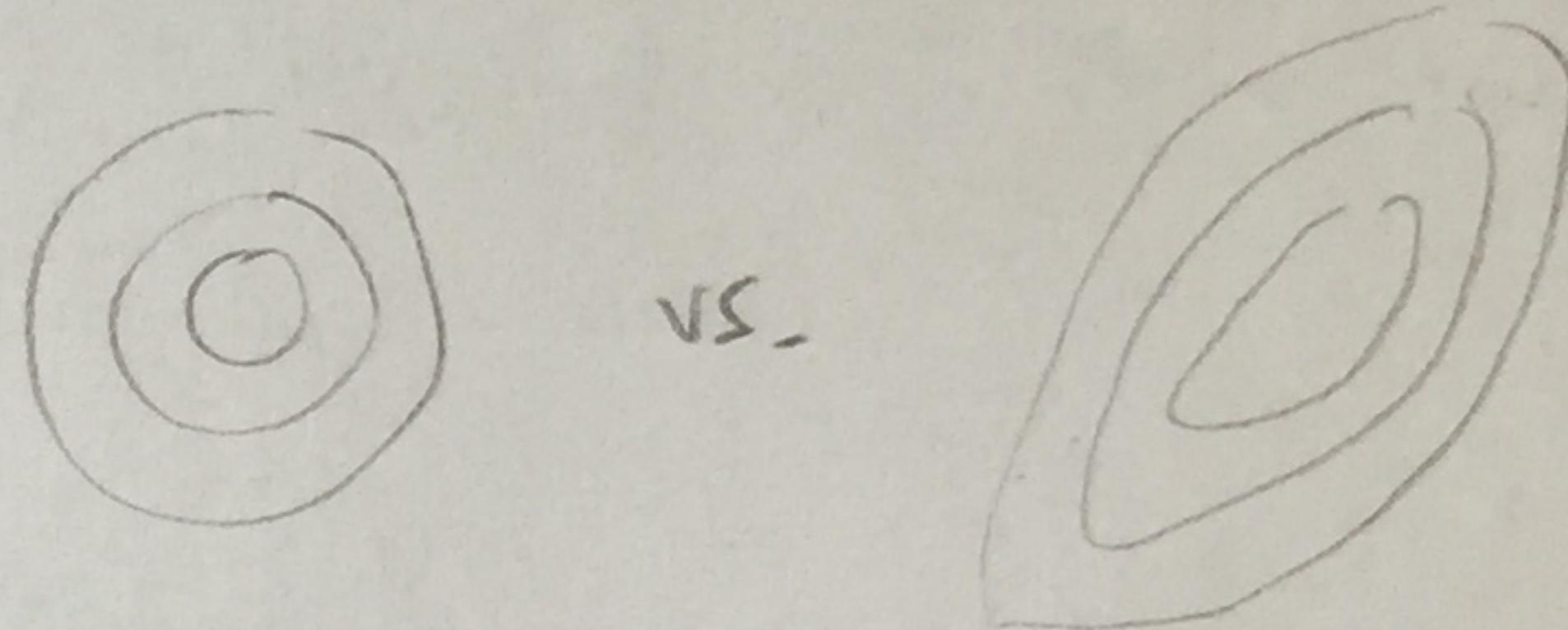
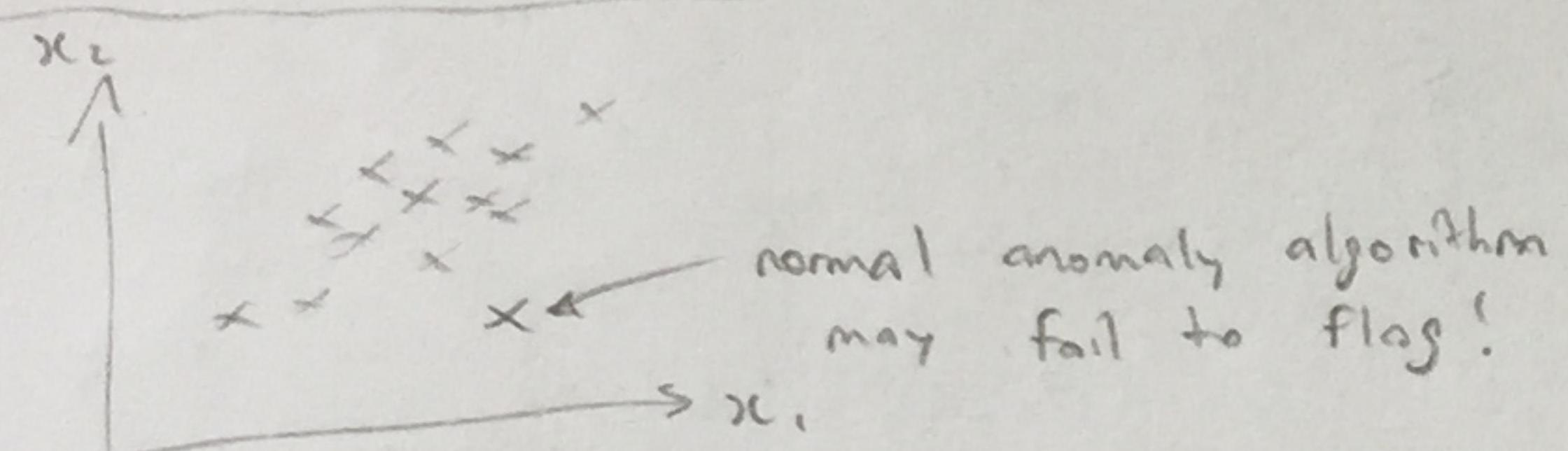
↳ most common problem: $p(x)_{\text{norm}} \approx p(x)_{\text{anom}}$

↳ can identify new/important features that solve this:

↳ choose features that might be large/small upon anomaly



Multivariate Gaussian Distributions:



⇒ $x \in \mathbb{R}^n \rightarrow$ model $p(\underline{x})$ all at once rather than $p(x_1), p(x_2)$, etc.

$$p(\underline{x}) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp \left(-\frac{1}{2} (\underline{x} - \underline{\mu})^T \Sigma^{-1} (\underline{x} - \underline{\mu}) \right) \quad \text{where } |\Sigma| = \text{determinant of } \Sigma \in \mathbb{R}^{n \times n}$$

$$\text{e.g. } \underline{\mu} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = \mathbb{I} \rightarrow \begin{array}{c} \text{contours of } p(\underline{x}) \\ \text{symmetric ellipses centered at } (0,0) \end{array}$$

$$\underline{\mu} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \Sigma = \begin{bmatrix} 0.6 & 0 \\ 0 & 0.6 \end{bmatrix} \rightarrow \begin{array}{c} \text{contours of } p(\underline{x}) \\ \text{ellipses centered at } (0,0) \end{array}$$

$$\underline{\mu} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \Sigma = \begin{bmatrix} 0.6 & 0 \\ 0 & 1 \end{bmatrix} \rightarrow \begin{array}{c} \text{contours of } p(\underline{x}) \\ \text{ellipses centered at } (0,0) \end{array}$$

⇒ x_2 has greater variance than x_1 .

$$\underline{\mu} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \Sigma = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix} \rightarrow \begin{array}{c} \text{contours of } p(\underline{x}) \\ \text{ellipses centered at } (0,0) \end{array}$$

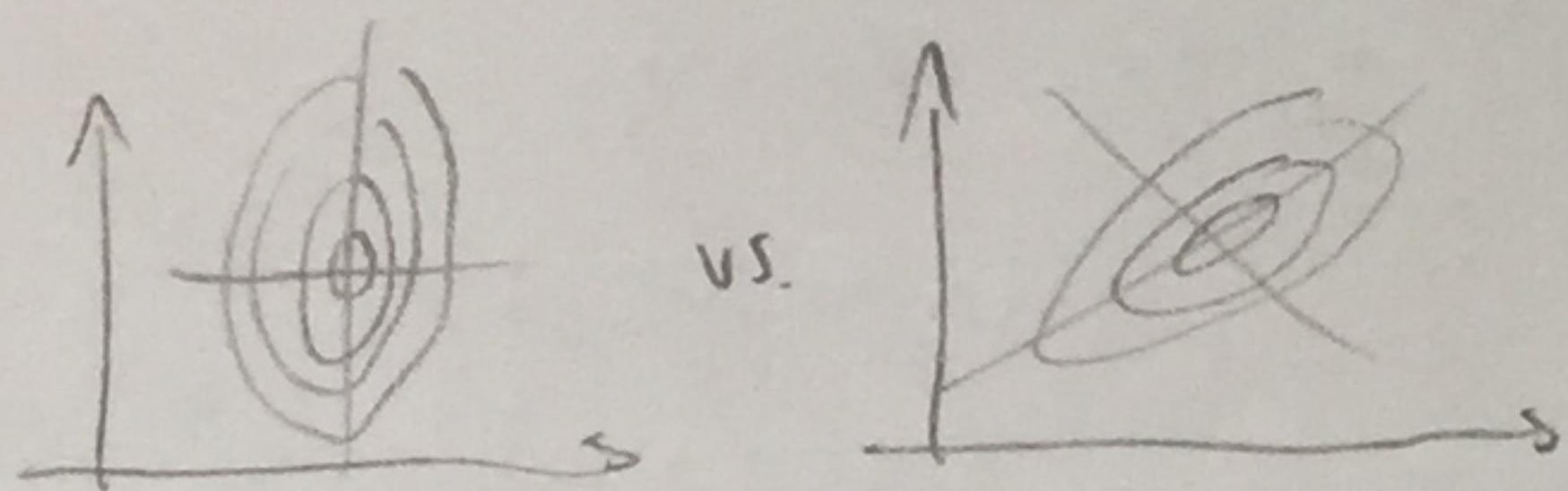
Note: $\underline{\mu}$ determines peak

Given $\underline{x} = [x^{(1)}, x^{(2)}, \dots, x^{(m)}]^T, x \in \mathbb{R}^n$:

$$\underline{\mu} = \frac{1}{m} \sum_{i=1}^m x^{(i)}, \quad \Sigma = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \underline{\mu})(x^{(i)} - \underline{\mu})^T$$

\Rightarrow Algorithm: get u & $\Sigma \rightarrow p(x_{\text{new}}) < \epsilon \rightarrow \text{anomaly!}$

\Rightarrow In original model $p(x) = \underbrace{p(x_1; m_1, b_1^2)}_{\text{each of these is multivariate special case.}} \times p(x_2; m_2, b_2^2) \dots$



AXIS ALIGNED! $\rightarrow \Sigma = \begin{bmatrix} a & 0 & 0 \\ 0 & b & 0 \\ 0 & 0 & c \end{bmatrix}$

\hookrightarrow multivariate automatically captures correlations between features!

\hookrightarrow computationally $\uparrow \downarrow \rightarrow \Sigma^{-1}$ where $\Sigma \in \mathbb{R}^{n \times n}$ (large n) $\Sigma \sim \frac{n^2}{2}$

\hookrightarrow also needs $m > n$, otherwise Σ non-invertible (want $m \geq 10n$)

Recommender Systems:

e.g. Predicting movie ratings: 0-5 stars
4 users (1-4), 5 movies

n_u : # users

n_m : # movies

$r(i,j) = 1$ if user j rated movie i

$y^{(i,j)}$ = rating by user j for movie i (if $r(i,j)=1$)

	user $\rightarrow n_u$			
rating	5	5	0	0
	↓			
n_m	5	?	?	0
	?	4	0	?
	0	0	5	4
	0	0	5	?

\Rightarrow want to fill in
? in order to
recommend movies!

\Rightarrow Have 2 features, x_1 (romance) & x_2 (action) for movies

\hookrightarrow movie 1: $x^{(1)} = [1 \ 0.9 \ 0]$ \rightarrow feature vectors for movie i

$$x^{(2)} = [1 \ 1.0 \ 0.01] \text{ etc.}$$

\Rightarrow for each user, learn $\theta^{(j)} \in \mathbb{R}^{n+1}$ ($= 3$ here) \rightarrow rating of user $j = (\theta^{(j)})^T x^{(i)}$
 \hookrightarrow parameter vector for user j

Formulation:

$$\min_{\theta^{(j)}} \frac{1}{2m^{(j)}} \sum_{i: r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2m^{(j)}} \sum_{k=1}^n (\theta_k^{(j)})^2$$

\hookrightarrow to learn all $\theta^{(j)}$'s:

$$\min_{\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(n_u)}} \frac{1}{2} \sum_{j=1}^{n_u} \sum_{i: r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^n (\theta_k^{(j)})^2$$

$J(\theta^{(1)}, \dots, \theta^{(n_u)})$

Gradient Descent for This:

$$\theta_k^{(j)} := \theta_k^{(j)} - \alpha \sum_{i: r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)}) x_k^{(i)} \quad \text{for } k=0$$

$$\theta_k^{(j)} := \theta_k^{(j)} - \alpha \left[\sum_{i: r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)}) x_k^{(i)} + \lambda \theta_k^{(j)} \right] \quad \text{for } k \neq 0$$

$\frac{\partial}{\partial \theta_k^{(j)}} J(\theta^{(0)}, \dots, \theta^{(n)})$

Collaborative Filtering: feature learning (self-learns features)

↳ where do x_1 (romance) & x_2 (action) come from?

↳ instead of having x_1, x_2 vals, survey users $\rightarrow \theta^{(j)}$ given by user

↳ given ratings of users & self-reported $\theta^{(j)}$'s, can classify movie type!

\Rightarrow given $\theta^{(0)}, \dots, \theta^{(n)}$, learn $x_1, \dots, x^{(nm)}$:

$$\min_{x^{(1)}, \dots, x^{(nm)}} \frac{1}{2} \sum_{i=1}^{nm} \sum_{j: r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{i=1}^{nm} \sum_{k=1}^n (x_k^{(i)})^2$$

so given x 's (ratings) \rightarrow predict θ 's (user preferences)

& give θ 's (user preferences) \rightarrow predict x 's (ratings)

↳ can guess $\theta \rightarrow$ get $x \rightarrow$ update $\theta \rightarrow$ get x etc.

\Rightarrow more efficient algorithm (simultaneous x & θ):

$$J(x^{(1)}, \dots, x^{(nm)}, \theta^{(0)}, \dots, \theta^{(n)}) = \frac{1}{2} \sum_{(i,j): r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{i=1}^{nm} \sum_{k=1}^n (x_k^{(i)})^2 + \frac{\lambda}{2} \sum_{j=1}^n \sum_{k=1}^n (\theta_k^{(j)})^2$$

↳ $\min_{x^{(1)}, \dots, x^{(nm)}, \theta^{(0)}, \dots, \theta^{(n)}} J(\dots)$

- j : users (n) $r(:,j)=1 \rightarrow$ rated / exists
- i : movies (nm)
- k : # genres (n) \rightarrow (features in θ)

① Initialize $x^{(1)}, \dots, x^{(nm)}, \theta^{(0)}, \dots, \theta^{(n)}$ to small random vals

② minimize $J(x, \theta) \rightarrow$ gradient descent: $x_k^{(i)} = x_k^{(i)} - \alpha \left[\sum_{j: r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)}) \theta_k^{(j)} + \lambda x_k^{(i)} \right]$

NOTE: no $x_0 = 1$, $x \in \mathbb{R}^n$
 $\theta \in \mathbb{R}^n$

$$\theta_k^{(j)} = \theta_k^{(j)} - \alpha \left[\sum_{i: r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)}) x_k^{(i)} + \lambda \theta_k^{(j)} \right]$$

- ③ For user w/ parameters θ & movie w/ learned features (genres) \underline{x} ,
predict rating of $\theta^T \underline{x} \rightarrow (\theta^{(j)})^T \underline{x}^{(i)}$

Vectorization of this Algorithm:

$$\begin{array}{l} \text{movie users} \rightarrow j \\ \downarrow i \quad n_u = 4 \end{array} \Rightarrow \underline{Y} = \begin{bmatrix} 5 & 5 & 0 & 0 \\ 5 & - & - & 0 \\ -4 & 0 & - & - \\ 0 & 0 & 5 & 4 \\ 0 & 0 & 5 & 0 \end{bmatrix} \in \mathbb{R}^{n_m \times n_u} \quad \left. \begin{array}{l} \text{y}^{(i,j)} \text{ entries} \\ \text{(actual)} \end{array} \right\}$$

\Rightarrow predicted rankings: $\downarrow \begin{bmatrix} (\theta^{(1)})^T \underline{x}^{(1)} & \dots & (\theta^{(n_u)})^T \underline{x}^{(1)} \\ \vdots & \ddots & \vdots \\ (\theta^{(1)})^T \underline{x}^{(n_m)} & \dots & (\theta^{(n_u)})^T \underline{x}^{(n_m)} \end{bmatrix} \rightarrow$ what user j predicts on movie i

$$\underline{X} = \begin{bmatrix} -(\underline{x}^{(1)})^T - \\ \vdots \\ -(\underline{x}^{(n_m)})^T - \end{bmatrix}, \quad \underline{\theta} = \begin{bmatrix} -(\theta^{(1)})^T - \\ \vdots \\ -(\theta^{(n_u)})^T - \end{bmatrix}$$

\hookrightarrow pred. rankings = $\underline{X} \underline{\theta}^T$ (Low-Rank Matrix Factorization)

\Rightarrow for each product (movie), we learn feature vector $\underline{x}^{(i)} \in \mathbb{R}^n$

\hookrightarrow e.g. $x_1 = \text{romance}$, $x_2 = \text{action}$, etc.

\hookrightarrow how to find movies j related to movie i ?

\hookrightarrow want $\|\underline{x}^{(i)} - \underline{x}^{(j)}\|$ small \rightarrow then movies i & j similar!

Mean Normalization:

$$\Rightarrow \text{If no ratings for user} \rightarrow \theta^{(\text{user})} = \underline{0}, \quad (\theta^{(\text{user})})^T \underline{x}^{(\text{movie})} = 0$$

\hookrightarrow to fix this, get average for each movie, then subtract this from all rated val -> avg. rating of each movie now 0

\Rightarrow Now, prediction = $(\theta^{(j)})^T \underline{x}^{(i)} + m_i \rightarrow$ for user w/ no ratings, give average rating m_i instead of 0

(If movie w/ no ratings, don't recommend it anyway)