

⇒ Week-9

⇒ Anomaly detection: broken motivation

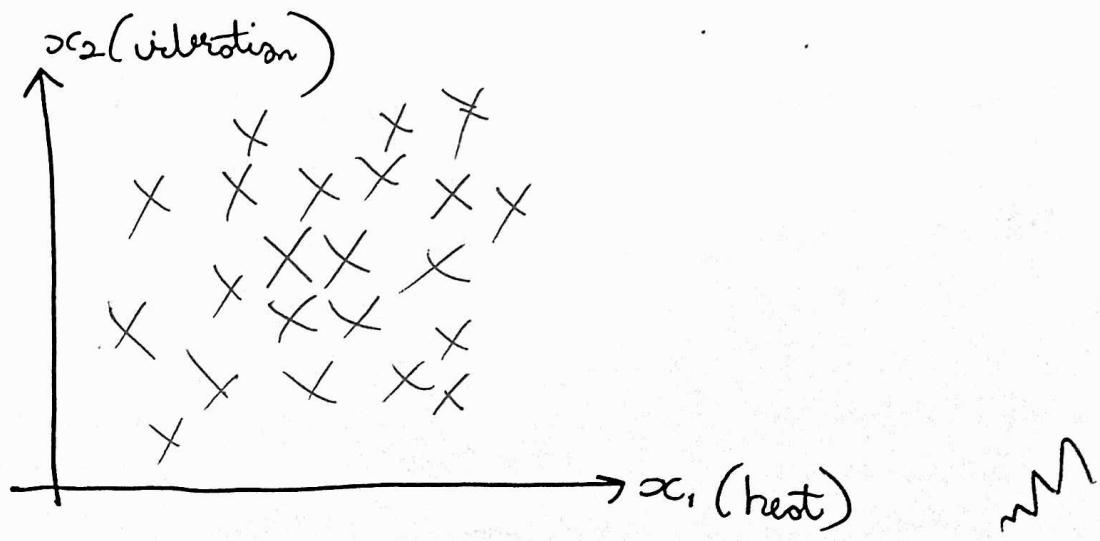
⇒ What is anomaly detection?

Anomaly detection refers to identification of items or events that do not conform to an expected pattern or to other items in a dataset.

Eg: Supp. we're on aircraft manufacturing company

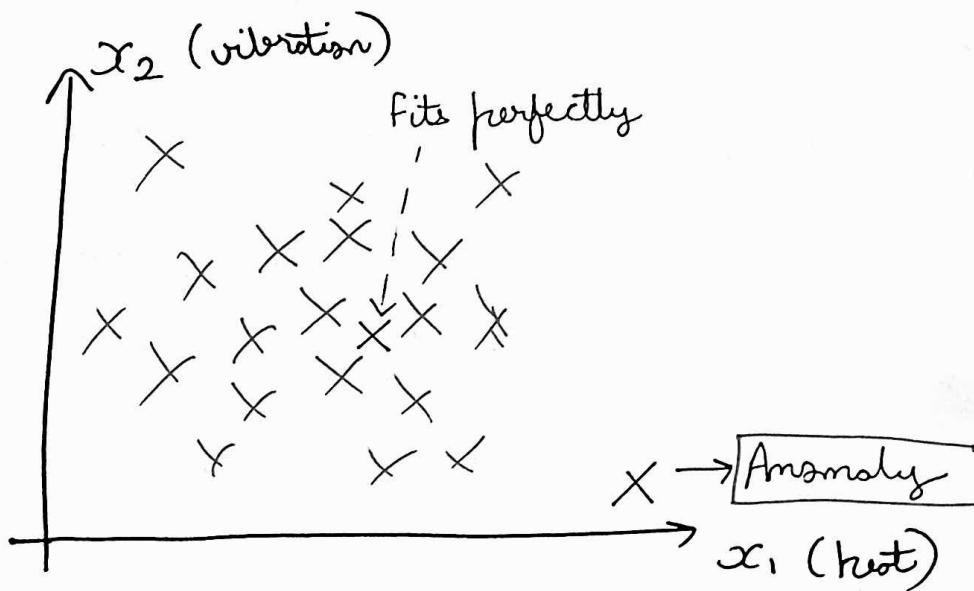
I have a dataset (x^1 to x^m) consisting of 2 features & examples: x_1 (heat) & x_2 (vibration)

Supp. our dataset looks like this:



→ Next day you have 2 new aircraft engines to test. When plotted the with the full dataset looks as follows:

(Note the new entries have been marked with green)



→ The anomaly detection method will see if the new engine is anomalous or not (when compared to the prev. engines)

for the above example, it'll

detect the anomalous aircraft & send it for further testing.

\Rightarrow Density estimation

\rightarrow In probability & statistics, density estimation is the construction of an estimate, based on observed data, of an unobservable underlying probability density function (PDF).

The unobservable PDF is thought of as the density according to which a large population is distributed.

\rightarrow Density estimation is also frequently used in anomaly detection or novelty ~~detection~~ detection. If an observation lies in a very low-density region, it is likely to be an anomaly or a novelty.

→ Given a dataset : $\{x^1, x^2, \dots, x^m\}$ s.t.

the dataset has non-anomalous data.

→ Given a new example x^{test} . We check if its anomalous or not.

→ How do we do it?

Using our training dataset, we build a model. We can access this model using $p(x)$.

* This $p(x)$ is not a probability but rather the normalized probability density function as parameterized by the feature vector x .

⇒ Determination of the actual probability can be done by integrating $p(x)$ over some range.

Q&A

* $p(x)$ indirectly asks "What is the probability that the example x is non-anomalous?"

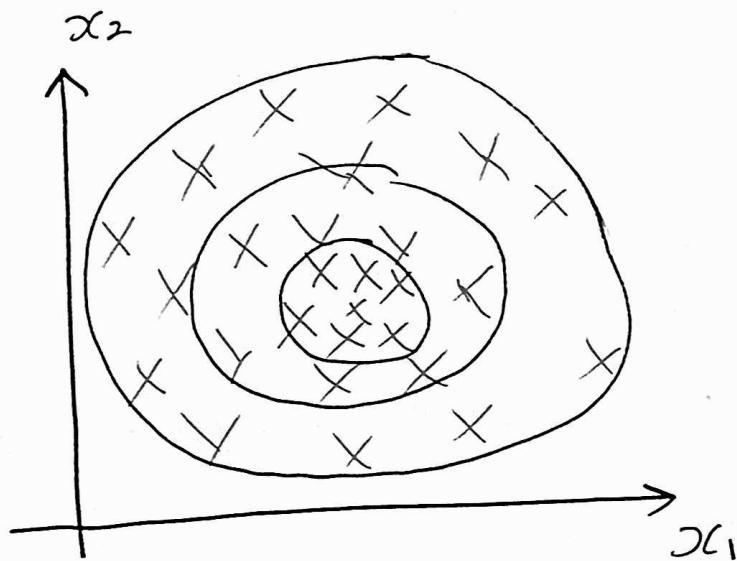
→ Once model is built,

* if $p(x^{\text{test}}) < \epsilon$ } flag this as an anomaly

* if $p(x^{\text{test}}) > \epsilon$ } flag this as non
anomalous

↳ ϵ is a threshold probability value
which we define depending on how
sure we need / want to be

→ For our aircraft example we want our
model to (graphically) look like this:



(The inner circles have more probability of getting
characterized as anomalies.)

⇒ Applications of anomaly detection

→ Fraud detection

* x^i = features of ' i 'th user's activities

eg of features : Typing speed, No. of webpage/transactions
Still time

* Model $\mu(x)$ from data

* Identify unusual users by checking which have $\mu(x) < \epsilon$.

→ Manufacturing : Discussed in aircraft example

→ Monitoring computers in data center

* $x^i \rightarrow$ Memory use, no. of diskaccesses/sec, CPU load,
Network traffic

* Identify crashing system by checking $\mu(x) < \epsilon$

→ If our anomaly detector is flagging too many anomalous examples, then we need to decrease our threshold ϵ .

\Rightarrow Gaussian distribution (or Normal distribution)

\rightarrow The Gaussian distribution is a familiar bell shaped curve that is described by the function

$$N(\mu, \sigma^2)$$

\rightarrow Formally,

Say $x \in \mathbb{R}^m$. If \tilde{x} $\xrightarrow{\text{dataset}}$ is a distributed Gaussian with mean μ , variance σ^2 , then:

$$\tilde{x} \sim N(\mu, \sigma^2)$$

* The little \sim or 'tilde' can be read as "distributed as".

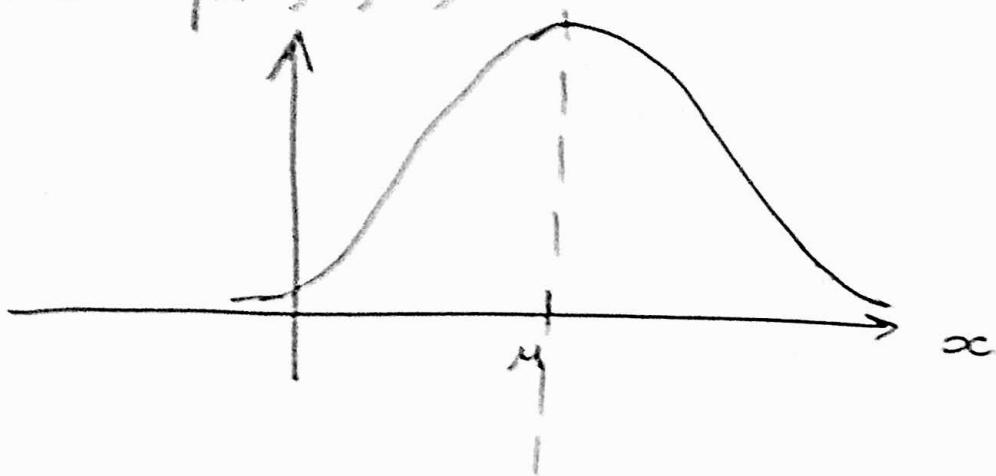
* The Gaussian distribution is completely parameterized by a mean (μ) \rightarrow describes the center of the curve & sigma standard deviation (σ) \rightarrow describes the width of the curve:

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

- Why standard deviation represents width?
- Standard deviation ⁽⁶⁾ represents is a measure of variability: How spread out are the values?
- While a measure of central tendency (like mean, median) describe the "typical" value of the dataset, measures of variability define how far away the data points tend to fall from the center.
- A low value of σ denotes data points being clustered tightly around the center (σ mean is in our case). A high value signifies that they tend to fall further away.
- By,*
- We can see how low σ causes less width of the distribution and a high σ causes the distribution to be more wide.

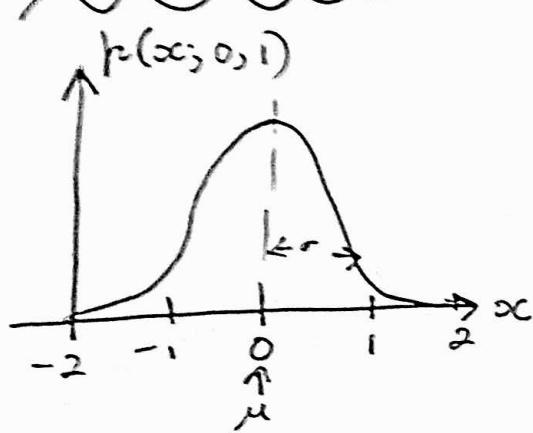
⇒ Gaussian distribution examples

→ Generalized $p(x; \mu, \sigma^2)$

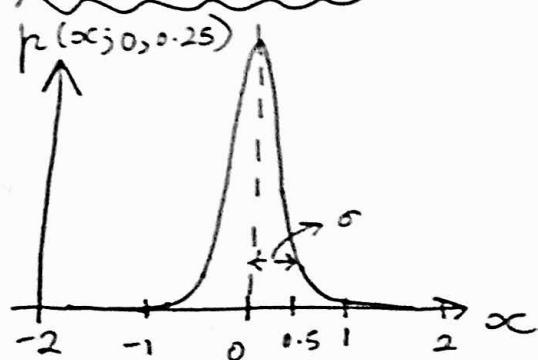


* Area is always equal to 1 but width is
center can be changed by playing with values of
 σ & μ respectively.

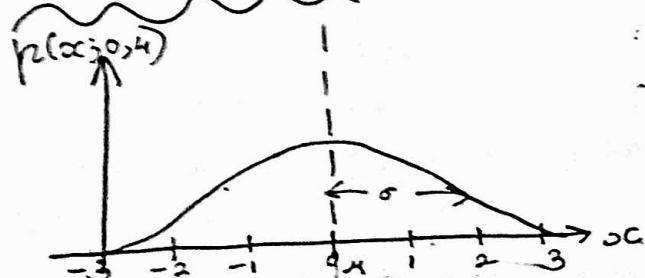
→ $\mu = 0, \sigma = 1$



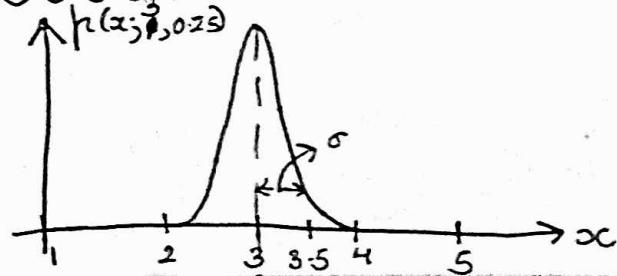
→ $\mu = 0, \sigma = 0.5$



→ $\mu = 0, \sigma = 2$



→ $\mu = 3, \sigma = 0.5$

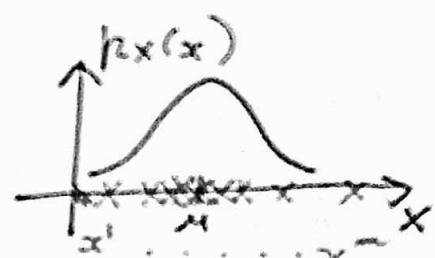


\Rightarrow Parameter estimation

- * Modelling our dataset in terms of a probabilistic model can help us in calculating probabilities for each example.
- * For anomaly detection, using Gaussian distribution makes sense ^{intuitively} because of its inherent structure & properties. eg: If you can visually tell the points which are closer to the "typical" or average value

\rightarrow Problem: Given a dataset $\{x^1, x^2, \dots, x^n\} \subset \mathbb{R}$
Can you estimate the distribution

\rightarrow Could be something like this:



- * Seems like a reasonable fit - data seems like a higher probability of being in the central region, lower probability of being further away

→ The Gaussian distribution can be completely described by μ & σ^2 as they can be calculated as follows:

$$\mu = \text{Average of examples} = \frac{1}{m} \sum_{i=1}^m x_i$$

$$\sigma^2 = \left(\text{Standard deviation} \right)^2 = \text{Variance} = \frac{1}{m} \sum_{i=1}^m (x_i - \mu)^2$$

* These values μ & σ^2 are

* These values are the maximum likelihood estimations (MLE) for μ & σ^2 .

* We can also do $\frac{1}{(m-1)}$ instead of $\frac{1}{m}$ which is used in statistics community, but in practice it makes hardly any difference.

- ⇒ Density estimation algorithm
- Before formally stating the anomaly detection algorithm, we need to estimate the probability density of our dataset.
- Given a training set $X : \{x^1, \dots, x^m\}$; ~~$x^i \in \mathbb{R}^n$~~
- where $x^i \in \mathbb{R}^n$
- * We model each of

* We model each of the features by assuming each feature is distributed according to a Gaussian distribution; s.t.

$$\boxed{\begin{aligned} X_1 &\sim N(\mu_1, \sigma_1^2) \\ X_2 &\sim N(\mu_2, \sigma_2^2) \\ &\vdots \\ X_n &\sim N(\mu_n, \sigma_n^2) \end{aligned}}$$

s.t. $p_{X_i}(x_i; \mu_i, \sigma_i^2) \rightarrow$ PDF of feature

x_i given μ_i & σ_i^2 using a Gaussian distribution.

→ We can assume that each feature is statistically independent from each other. (If they aren't, just use a dimensionality reduction technique like PCA to get independent features)

→ Since our features are independent, we can formulate the joint PDF of these features as the product of the individual PDFs.

~~Defn of PDF~~

$$p_{x_1, \dots, x_m}(x_1, \dots, x_m) = p_{x_1}(x_1; \mu_1, \sigma_1^2) \cdots p_{x_m}(x_m; \mu_m, \sigma_m^2)$$

→ ~~And we formulated the above~~

$$\therefore p_{x_1, \dots, x_m}(x_1, \dots, x_m) = \prod_{j=1}^m p_{x_j}(x_j; \mu_j, \sigma_j^2)$$

⇒ Anomaly detection algorithms

- 1) Choose features x_i that you think might be indicative of anomalous examples.
- 2) Fit (or calculate) parameters $\mu_1, \dots, \mu_m, \sigma_1^2, \dots, \sigma_m^2$

$$\mu_g = \frac{1}{m} \sum_{i=1}^m x_g^i$$
$$\sigma_g^2 = \frac{1}{m} \sum_{i=1}^m (x_g^i - \mu_g)^2$$

$$\rightarrow \text{Can represent } \boldsymbol{\mu} = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_m \end{bmatrix} = \frac{1}{m} \sum_{i=1}^m \cancel{x^i}$$

- 3) Given new example x , calculate $p(x)$

$$p(x) \equiv p_{x_1, \dots, x_m}(x_1, \dots, x_m) = \prod_{g=1}^m p_{x_g}(x_g; \mu_g, \sigma_g^2)$$
$$= \prod_{g=1}^m \frac{1}{\sqrt{2\pi} \sigma_g} \exp\left(-\frac{(x_g - \mu_g)^2}{2\sigma_g^2}\right)$$

Detect anomaly if $p(x) < \epsilon$

\Rightarrow Developing & evaluating an anomaly detection system

\rightarrow When developing a learning algo, making decisions (like crossing ^{which to include} "features", etc) becomes easier if we have a way of evaluating our learning algorithm.

\rightarrow In case of anomaly detection, we take some labelled data, categorized into ~~as~~ anomalous & non-anomalous examples ($y=0$ if ~~is~~ non-anomalous, $y=1$ if anomalous)

* Among that, take a large proportion of good, non-anomalous data for the training set on which to train $f(x)$.

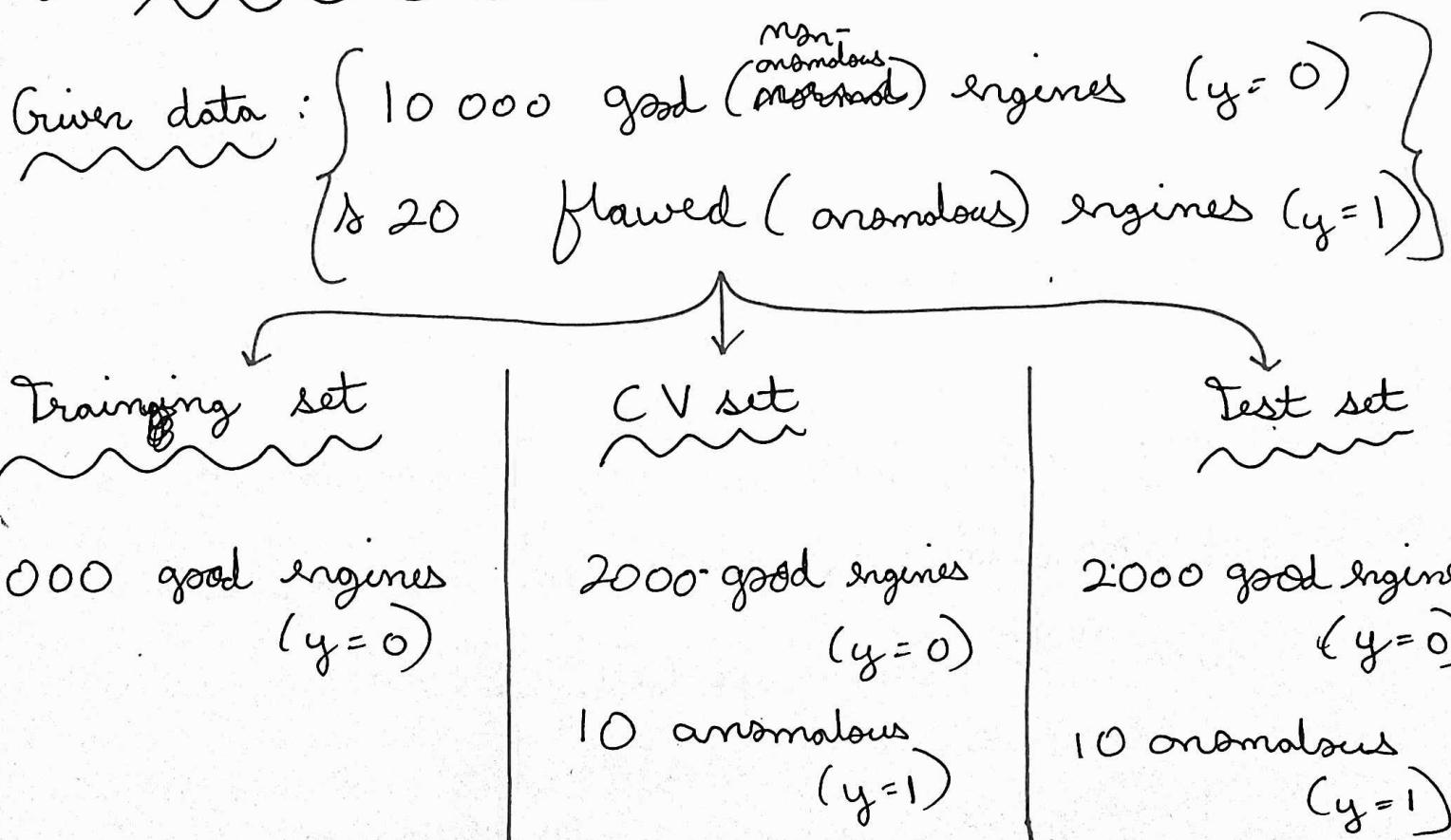
Training set : x^1, x^2, \dots, x^m (assume non anomalous)

→ Then take a smaller proportion of mixed anomalous & non-anomalous examples for your cross validation & test sets.

Cross validation set: $(x_{cv}^1, y_{cv}^1), \dots, (x_{cv}^{m_{cv}}, y_{cv}^{m_{cv}})$

Test set: $(x_{test}^1, y_{test}^1), \dots, (x_{test}^{m_{test}}, y_{test}^{m_{test}})$

→ eg: Aircraft engines example



→ In Summary:

- * Split the data 60 : 20 : 20 :: Training : CV : test set set set
- * Split the anomalous examples 50 : 50 :: CV : test set set

⇒ Algorithm evaluation workflow

- 1) Fit model $p_x(x)$ on training set $X = \{x^1, \dots, x^m\}$
- 2) On a cross validation/test set example x , predict

$$y = \begin{cases} 1, & p_x(x) < \epsilon \Rightarrow \text{anomaly} \\ 0, & p_x(x) \geq \epsilon \Rightarrow \text{non-anomalous} \end{cases}$$

- 3) Possible evaluation metrics: Since the no. of non-anomalous examples will always be much higher than the anomalous ones (or the classes are skewed), the alg that predicts $y=0$ will always have higher accuracy if classification accuracy is used. The foll. metrics will give better results:
 - * True+ve, False+ve, True-ve, False-ve
 - * Precision/Recall
 - * F₁ Score

→ For taking important decisions like choosing ϵ , which features to drop, use the cross validation set for that & choose the ~~configur~~ configuration which minimizes the cross validation error.

⇒ Anomaly detection vs Supervised learning

* Since, we've have or labelled data with ourselves, with examples that we know are either anomalous or non-anomalous, why don't we use a supervised learning technique like SVM, logistic regression or neural nets to directly learn from our data in predicting whether $Y=0$ or $Y=1$?

Anomaly detection

→ Use this when you have a very small no. of +ve examples ($y=1$)
(0-20 is common)

→ Reason: There can exist "infinitely" many different "types" of anomalies. Our low no. of labels for anomalous examples in our dataset doesn't help either!

It'd be hard for a supervised algo to predict a label for some anomaly because a future anomaly might look completely diff. than our example anomalies

Supervised learning

→ Use this when you have large no. of both +ve & -ve examples

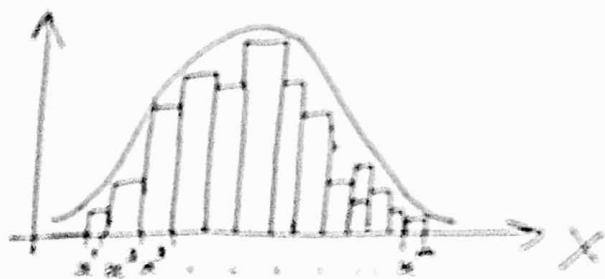
→ Reason: We have enough examples +ve (or anomaly) & examples for our algo to generalize for the every-type of +ve examples thrown at it in the future.

→ Uses of supervised learning: ~~dates~~ Email

Spam classification, Weather prediction (sunny/rainy/etc), Cancer classification

→ Uses of anomaly detection: Fraud detection, Manufacturing (e.g. Aircraft engines), Monitoring machines

- ⇒ Choosing what features to use
- The features you choose to use/include in your model will greatly impact how well your anomaly detection algorithm works.
- * We can check whether our features are "Gaussian" by plotting a histogram of our data & check for the bell shaped curve.



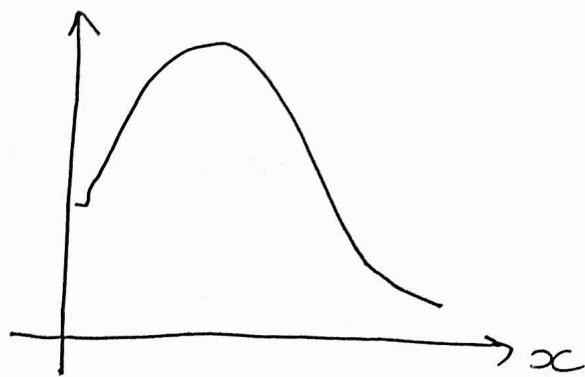
* You can decrease the no. of bins for a more "finer" visualization.

\Rightarrow Non Gaussian features

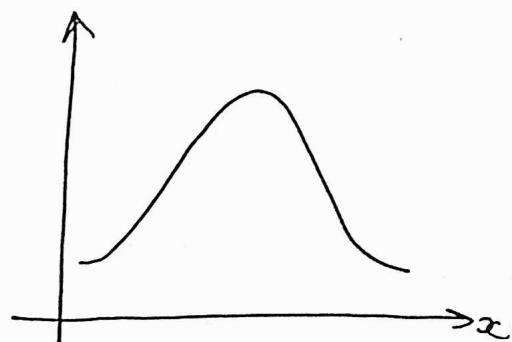
\Rightarrow We can apply some "transforms" on a example feature x that does not have the bell shaped curve:

- * $\log(x + c)$; $c \rightarrow \text{constant}$
- * \sqrt{x}
- * $x^{1/3}$

e.g.: For a left skewed distribution:



$\log(x)$



- ⇒ Error analysis for anomaly detection
- Requirement:
- * Want $p(x)$ to be large for non-anomalous examples x
 - * " " " " small " anomalous examples x
- A common problem:
- After training $p(x)$ might come out to be comparable relatively similar both between anomalous & non-anomalous examples alike.
- ☞ (Very similar to what happens in underfitting!)
- * In this case, choose examine the anomalous examples x that are giving a "higher" probability & try to figure out new features that will better distinguish the data.

→ eg.: Monitoring computers in a data center

Supp. our features are as follows:

$$x_1 = \text{memory use}$$

$$x_2 = \text{no. of (disc access)/sec}$$

$$x_3 = \text{CPU load}$$

$$x_4 = \text{network traffic}$$

Say, our server is stuck in an infinite loop, so our CPU load will shoot up but our network traffic is still low.

* Therefore $p(x_3) \rightarrow$ denoting CPU load decreases, but $p(x_4) \rightarrow$ denoting network traffic is at a higher value.

* As for calculating the final probability we'll multiply each probability ($p(x_1, x_2, x_3, x_4) = p(x_1)p(x_2)p(x_3)p(x_4)$)

So the higher val. of $p(x_4)$ brings down the lower val. of $p(x_3)$.

thus misclassifying the example as non-anomalous!

- There exists 2 sol's for this:
- 1) Choose a feature that highlights the unusually high CPU load & address the low network traffic at the same time.
eg: $x_5 = \frac{\text{CPU load}}{\text{Network traffic}}$; $x_6 = \frac{(\text{CPU load})^3}{\text{Network traffic}}$
 - 2) Use a multivariate Gaussian distribution [Coming up!]
- In summary, choose / create features that might take on unusually large or small values in the event of an anomaly.