

Day - 95, March 5, 2025 (Folgen 21, 2021)

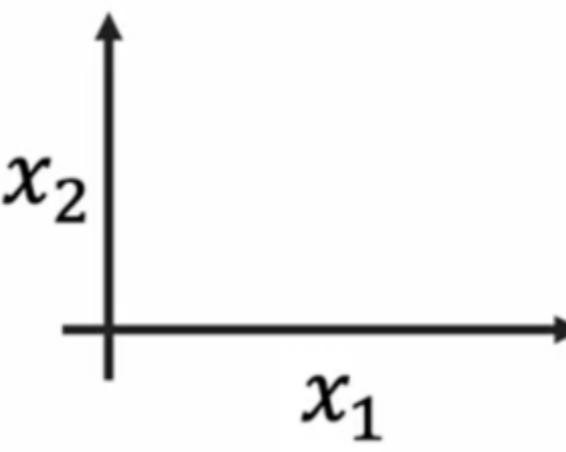
- ① Anomaly Detection vs. Supervised learning
- ② Choosing what features to use
- ③ Making Recommendations
- ④ Anomaly Detection Algorithm
- ⑤ Developing and Evaluating an Anomaly Detection Algorithm -

Source: Unsupervised Learning, Recommendation, Reinforcement Learning
Coursera.

Density estimation

Training set: $\{\vec{x}^{(1)}, \vec{x}^{(2)}, \dots, \vec{x}^{(m)}\}$

Each example $\vec{x}^{(i)}$ has n features


$$\vec{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

Assuming independent variables.

$$p(\vec{x}) = p(x_1; \mu_1, \sigma_1^2) * p(x_2; \mu_2, \sigma_2^2) * p(x_3; \mu_3, \sigma_3^2) * \dots * p(x_n; \mu_n, \sigma_n^2)$$

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

\sum
"odd"
~~Multiply~~

$$p(x_1 = \text{high temp}) = 1/10$$

$$p(x_2 = \text{high Vibra}) = 1/20$$

$$p(x_1, x_2) = p(x_1) * p(x_2)$$

$$= \frac{1}{10} * \frac{1}{20} \Rightarrow \frac{1}{200}$$

Anomaly detection algorithm

1. Choose n features x_i that you think might be indicative of anomalous examples.
2. Fit parameters $\mu_1, \dots, \mu_n, \sigma_1^2, \dots, \sigma_n^2$

$$\bar{x}_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)}$$

$$\sigma_j^2 = \frac{1}{m} \sum_{i=1}^m (x_j^{(i)} - \bar{x}_j)^2$$

vectorized formula.

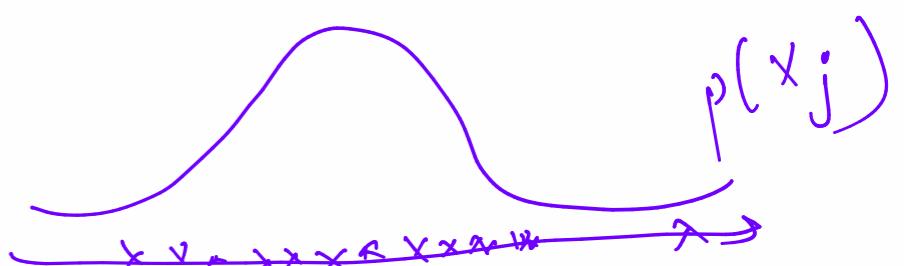
$$\bar{x} = \frac{1}{m} \sum_{i=1}^m x^{(i)}$$

$$\bar{x} = \begin{bmatrix} \bar{x}_1 \\ \bar{x}_2 \\ \vdots \\ \bar{x}_n \end{bmatrix}$$

3. Given new example x , compute $p(x)$:

$$p(x) = \prod_{j=1}^n p(x_j; \bar{x}_j, \sigma_j^2)$$
$$= \prod_{j=1}^n \exp \left(-\frac{(x_j - \bar{x}_j)^2}{2\sigma_j^2} \right)$$

→ Anomaly if $p(x) < \epsilon$



Step 5:
 $p(\vec{x})$ = Independent probability.

- ① Feature Selection (x_1, x_2, \dots, x_n)
- ② Model Assumption (Normal dist.)

Density estimation

Training set: $\{\vec{x}^{(1)}, \vec{x}^{(2)}, \dots, \vec{x}^{(m)}\}$

Each example $\vec{x}^{(i)}$ has n features

$$\vec{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

$$p(\vec{x}) = p(x_1; \mu_1, \sigma_1^2) * p(x_2; \mu_2, \sigma_2^2) * p(x_3; \mu_3, \sigma_3^2) * \dots * p(x_n; \mu_n, \sigma_n^2)$$

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

\sum "add" \prod "multiply"

$$p(x_1 = \text{high temp}) = 1/10$$

$$p(x_2 = \text{high vibra}) = 1/20$$

$$p(x_1, x_2) = p(x_1) * p(x_2)$$

$$= \frac{1}{10} * \frac{1}{20} = \frac{1}{200}$$

Anomaly detection algorithm

1. Choose n features x_i that you think might be indicative of anomalous examples.

2. Fit parameters $\mu_1, \dots, \mu_n, \sigma_1^2, \dots, \sigma_n^2$

$$\mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)}$$

$$\sigma_j^2 = \frac{1}{m} \sum_{i=1}^m (x_j^{(i)} - \mu_j)^2$$

3. Given new example x , compute $p(x)$:

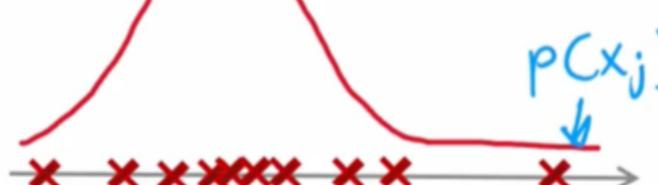
$$p(x) = \prod_{j=1}^n p(x_j; \mu_j, \sigma_j^2) = \prod_{j=1}^n \frac{1}{\sqrt{2\pi}\sigma_j} \exp\left(-\frac{(x_j - \mu_j)^2}{2\sigma_j^2}\right)$$

Anomaly if $p(x) < \epsilon$

Vectorized formula

$$\vec{\mu} = \frac{1}{m} \sum_{i=1}^m \vec{x}^{(i)}$$

$$\vec{\mu} = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{bmatrix}$$



Andrew Ng

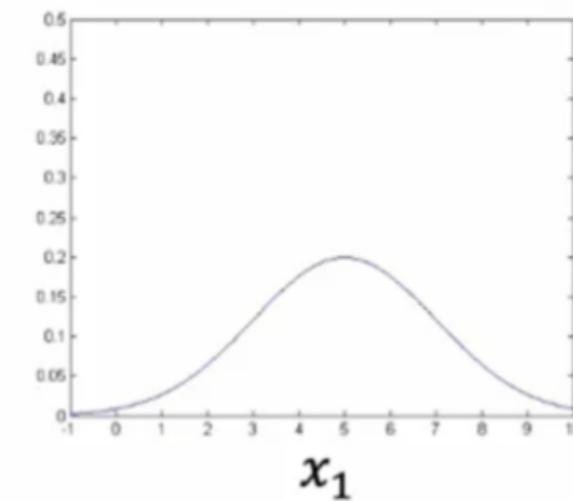
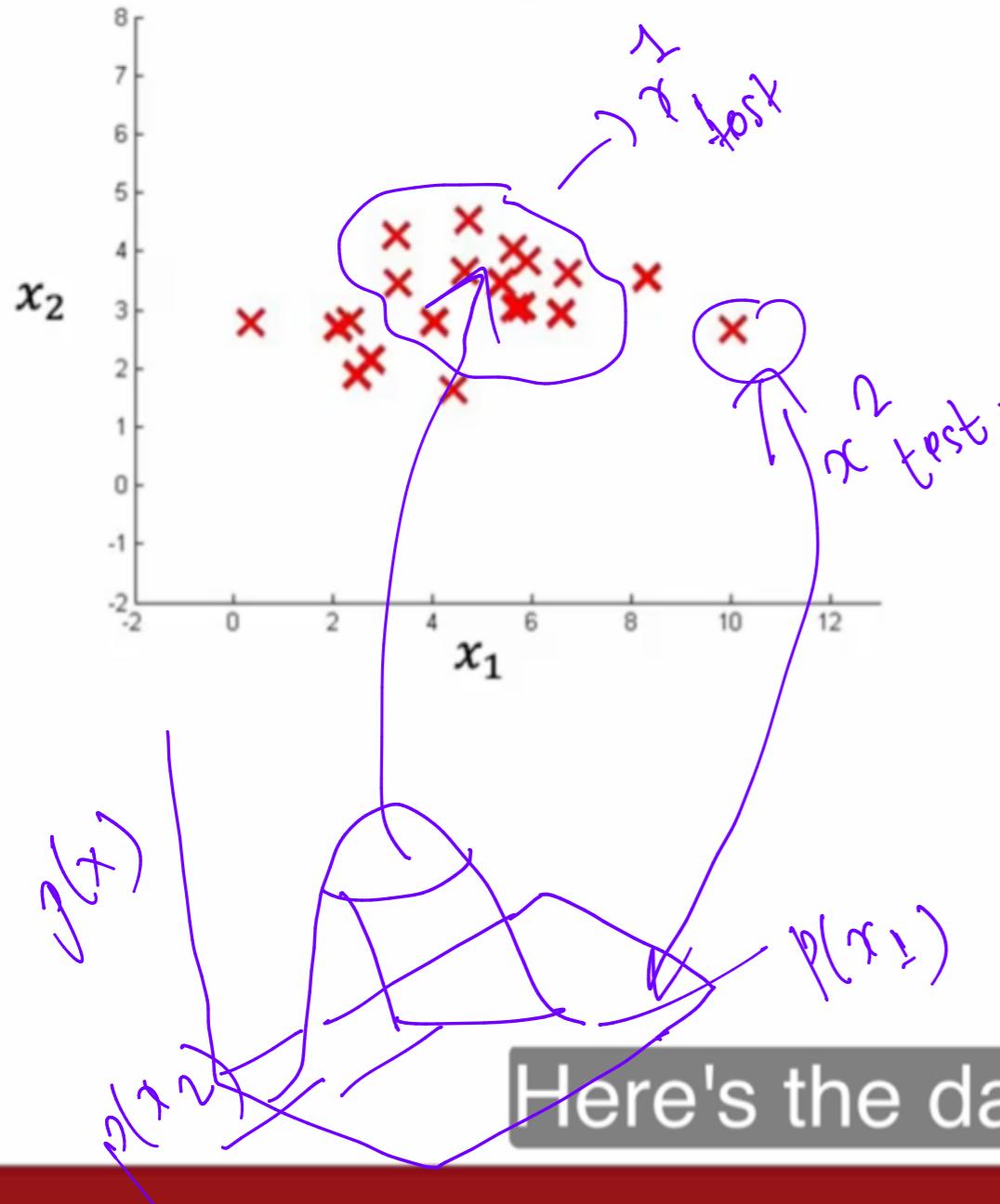
③ Parameter Estimation
 μ_j and σ_j^2

④ Probability Calculation:
 $p(x) = \prod_{j=1}^n \frac{1}{\sqrt{2\pi}\sigma_j} \exp\left(-\frac{(x_j - \mu_j)^2}{2\sigma_j^2}\right)$

⑤ Anomaly Detection

⑥ Intuition: See distribution.

Anomaly detection example



$$p(x_1; \mu_1, \sigma_1^2)$$

$$\mu_1 = 5, \sigma_1 = 2$$

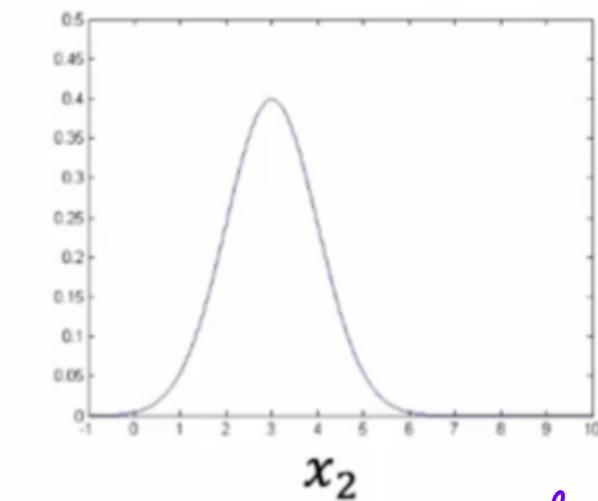
$$P(x^{(2)}_{\text{test}}) = 0.0426$$

"ok"

↓

Not ok:

x_2



$$p(x_2; \mu_2, \sigma_2^2)$$

$$\mu_2 = 3, \sigma_2 = 2$$

$$P(x^{(2)}_{\text{test}}) = 0.0021$$

↓

Not ok:

Anomaly detection example

Summary of Gaussian Anomaly Detection Algorithm

- Feature Selection:** Choose features x_1, x_2, \dots, x_n that may indicate anomalies (e.g., engine heat and vibrations).
- Density Estimation:** Model the probability $p(x)$ of a feature vector x using the product of individual feature probabilities:

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

This assumes feature independence, but the algorithm works even if this assumption is not strictly true.

- Parameter Estimation:** Compute mean μ_j and variance σ_j^2 for each feature x_j using training data:

$$\mu_j = \frac{1}{m} \sum_{i=1}^m x_{ij}, \quad \sigma_j^2 = \frac{1}{m} \sum_{i=1}^m (x_{ij} - \mu_j)^2$$

- Probability Calculation:** Compute the Gaussian probability density function (PDF) for each feature:

$$p(x_j) = \frac{1}{\sqrt{2\pi\sigma_j^2}} \exp\left(-\frac{(x_j - \mu_j)^2}{2\sigma_j^2}\right)$$

Then compute $p(x)$ by multiplying probabilities for all features.

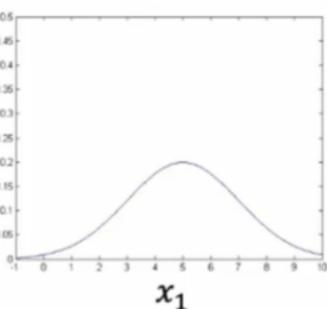
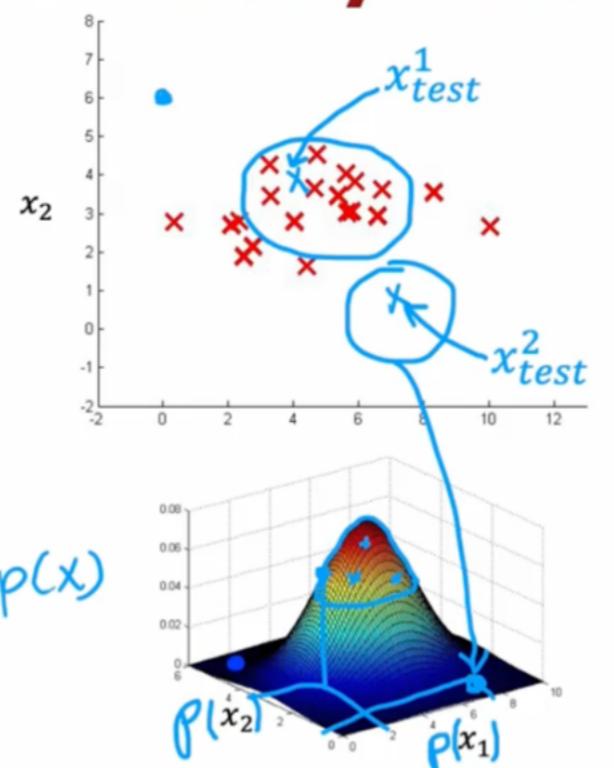
- Anomaly Detection:** Compare $p(x)$ with a threshold ϵ :

- If $p(x) < \epsilon$, flag x as an anomaly.
- If $p(x) \geq \epsilon$, classify x as normal.

- Intuition:** If any feature value is significantly different from the training distribution, $p(x)$ becomes very small, indicating an anomaly.

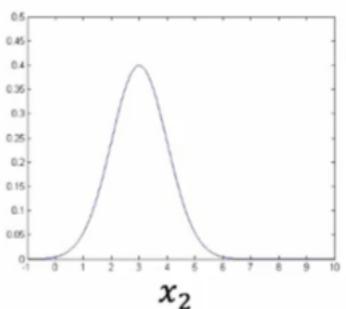
7. Visualization:

→ 2D Gaussian distributions show high-probability regions in the center from many data centers.



$$\mu_1 = 5, \sigma_1 = 2$$

$p(x_1; \mu_1, \sigma_1^2)$



$$\mu_2 = 3, \sigma_2 = 1$$

$p(x_2; \mu_2, \sigma_2^2)$

$$\epsilon = 0.02$$

$$p(x_{test}^{(1)}) = 0.0426 \rightarrow "ok"$$

$$p(x_{test}^{(2)}) = 0.0021 \rightarrow \text{anomaly}$$

DeepLearning.AI Stanford ONLINE

Andrew Ng

~~PDF~~ $p(x) = \prod_{j=1}^n p(x_j | \mu_j, \sigma_j^2)$

$$p(x_j) = \frac{1}{\sqrt{2\pi\sigma_j^2}} e^{-\frac{(x_j - \mu_j)^2}{2\sigma_j^2}}$$

The importance of real-number evaluation

When developing a learning algorithm (choosing features, etc.), making decisions is much easier if we have a way of evaluating our learning algorithm.

Assume we have some labeled data, of anomalous and non-anomalous examples.

Training set: $x^{(1)}, x^{(2)}, \dots, x^{(m)}$ ($y = 1$ for all normal examples; $y = 0$ for all anomalous examples)

Cross Validation set: $(x_{cv}^{(1)}, y_{cv}^{(1)}) \dots (x_{cv}^{(m_{cv})}, y_{cv}^{(m_{cv})})$ } include a few anomalous examples
Test set: $(x_{test}^{(1)}, y_{test}^{(1)}) \dots (x_{test}^{(m_{test})}, y_{test}^{(m_{test})})$ } mostly normal examples
 $y = 1$

Aircraft engines monitoring example

10000 good (normal) engines
20 flawed engines (anomalous)

Training set: 6000 $y=0$ good engines

(v: 2000 GE ($y=0$))

Test: 2000 GE ($y=0$), 10 anomalous ($y=1$)

Alternative: No test set

Training Set: 6000 good engines

(v: 4000 good engines ($y=0$), 20 anomalous ($y=1$))

tune ϵ

tune x_j

$i_2 \text{ to } 50$

$y=1$

train algorithm on training set.

Cross Validation set tune ϵ

(Using loss flawed engine is beneficial)

→ Smaller Data-Sets are prone to the higher risk of overfitting

The importance of real-number evaluation

When developing a learning algorithm (choosing features, etc.), making decisions is much easier if we have a way of evaluating our learning algorithm.

Assume we have some labeled data, of anomalous and non-anomalous examples.

$$y=1 \quad y=0$$

Training set: $x^{(1)}, x^{(2)}, \dots, x^{(m)}$ (assume normal examples/not anomalous)

$y=0$ for all training examples

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

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

} include a few anomalous examples
mostly normal examples
 $y=1$ $y=0$

Aircraft engines monitoring example

10000 good (normal) engines
20 flawed engines (anomalous)
2

Training set: 6000 good engines

2 to 50
 $y=1$

train algorithm on training set

CV: 2000 good engines ($y = 0$)
use cross validation set

Test: 2000 good engines ($y = 0$),

10 anomalous ($y = 1$)
tune ϵ tune x_j
10 anomalous ($y = 1$)

Alternative: No test set Use if very few labeled anomalous examples

Training set: 6000 good engines 2 higher risk of overfitting

CV: 4000 good engines ($y = 0$), 20 anomalous ($y = 1$)

tune ϵ tune x_j

Real-Number Evaluation:

↳ Quickly Assess and tune parameters like Epsilon

Labeled & Anomalies for Validation:

↳ Use a small set of known anomalies to define the model.

Data Partitioning:

↳ Training: Mostly Normal Data

↳ Cross-Validation: Mix of Normal and Anomalies for Tuning.

Test: Similar to Cross-Validation

for final evaluation.

↳ Merge Cross Validation and Test sets but be aware of Overfit.

Algorithm evaluation

Fit model $p(x)$ on training set $x^{(1)}, x^{(2)}, \dots, x^{(m)}$

On a cross validation/test example x , predict

$$y = \begin{cases} 1 & \text{if } p(x) < \varepsilon \text{ (anomaly)} \\ 0 & \text{if } p(x) \geq \varepsilon \text{ (normal)} \end{cases}$$

10
2000

Possible evaluation metrics:

- TP | FP | FN | TN

- Precision | Recall

→ F_1 - Score

Use Cross-Validation Set to choose parameter ε .

More Uncertainties

Anomaly detection

vs. Supervised learning

Very small number of positive examples ($y = 1$). (0-20 is common).

Large number of negative ($y = 0$) examples.

$p(x)$

$y=1$

Many different “types” of anomalies. Hard for any algorithm to learn from positive examples what the anomalies look like; future anomalies may look nothing like any of the anomalous examples we've seen so far.

Fraud:

Many diff forms
typos / Scenarios
whose

engine
might go
wrong

Large number of positive and negative examples.

20 positive examples

Enough positive examples for algorithm to get a sense of what positive examples are like, future positive examples likely to be similar to ones in training set.

Example: Spam

Anomaly detection

- Fraud detection
- Manufacturing - Finding new previously unseen defects in manufacturing. (e.g. aircraft engines)
→ Even brand new aircraft.
- Monitoring machines in a data center

Brand new
Known Examples
⋮
Not Known Examples

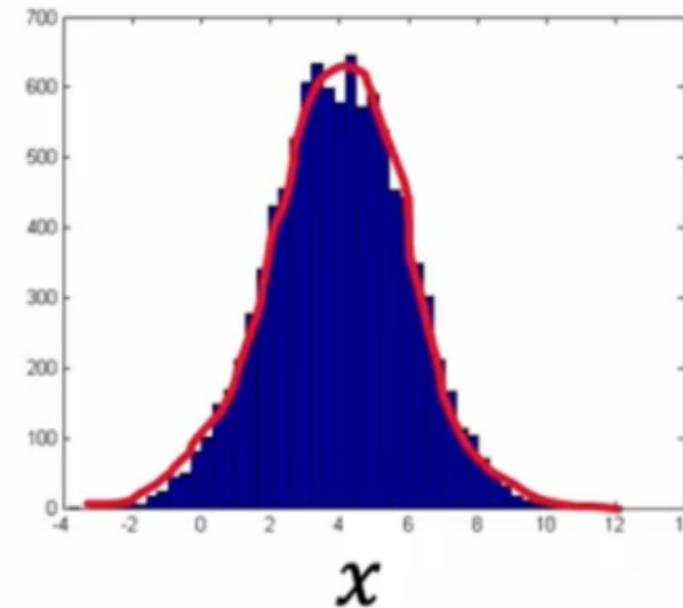
vs. Supervised learning

- Email spam classification
- Manufacturing - Finding known, previously seen defects
→ Mobile Bot Scratches Smartphone
- Weather prediction (sunny/rainy/etc.)
- Diseases classification

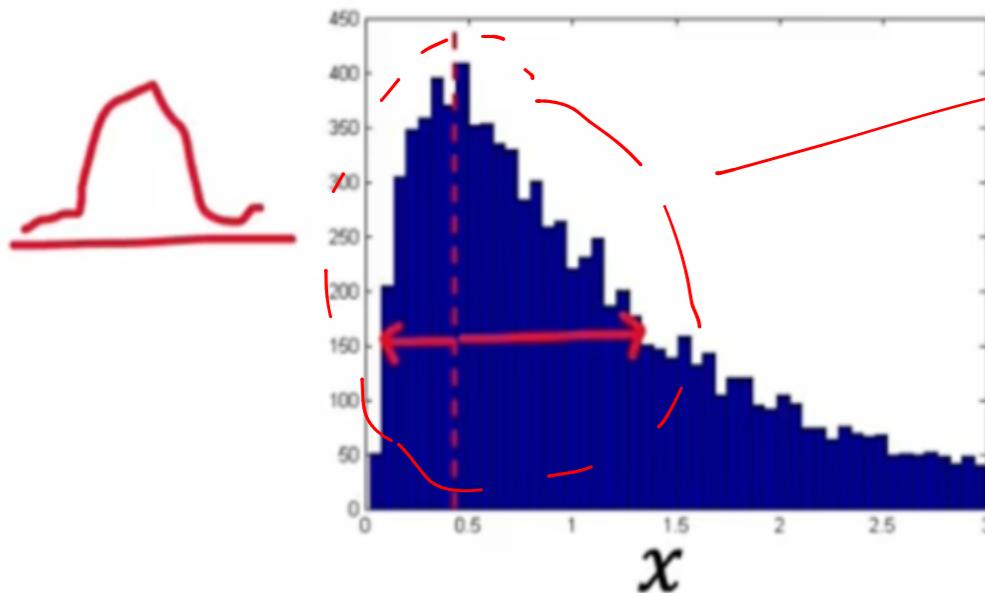
⋮
Known Examples

Always use Gaussian features (Map features into Normal distribution).

Non-gaussian features

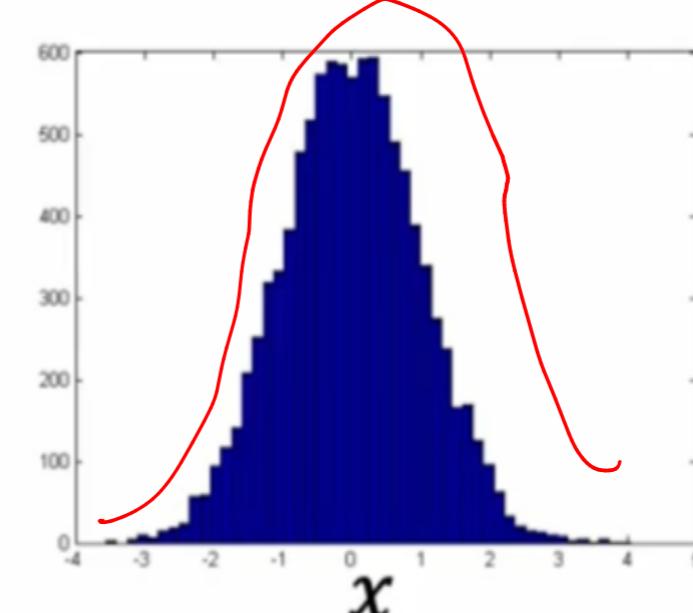


$P(x_1; \mu_1, \sigma_1^2)$
plt.hist(x)



use this.
np.log(x)

$$\begin{aligned}x_1 &\leftarrow \log(x_1) \\x_2 &\leftarrow \log(x_2 + 1) \\x_3 &= \sqrt{x_3} \rightarrow x_3^{1/2} \\x_4 &\leftarrow x_4^u\end{aligned}$$



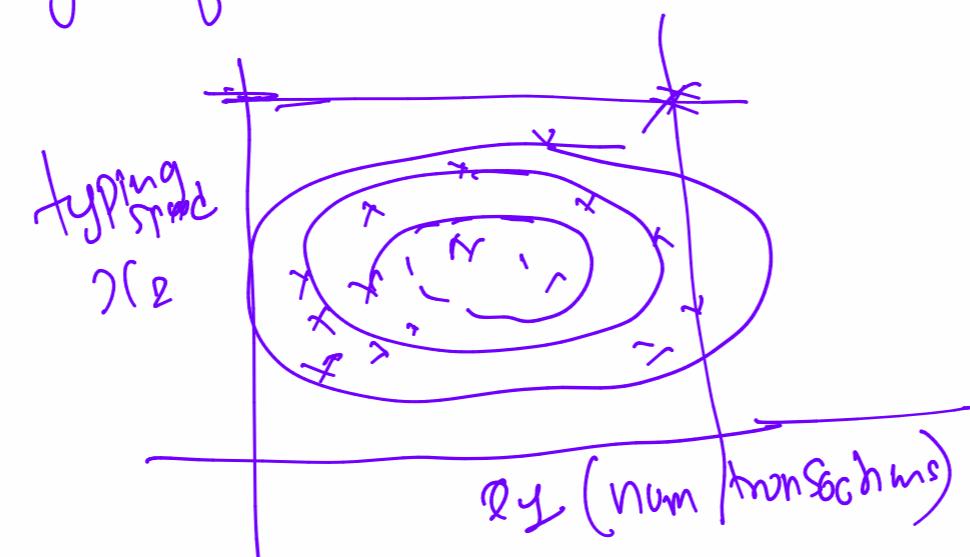
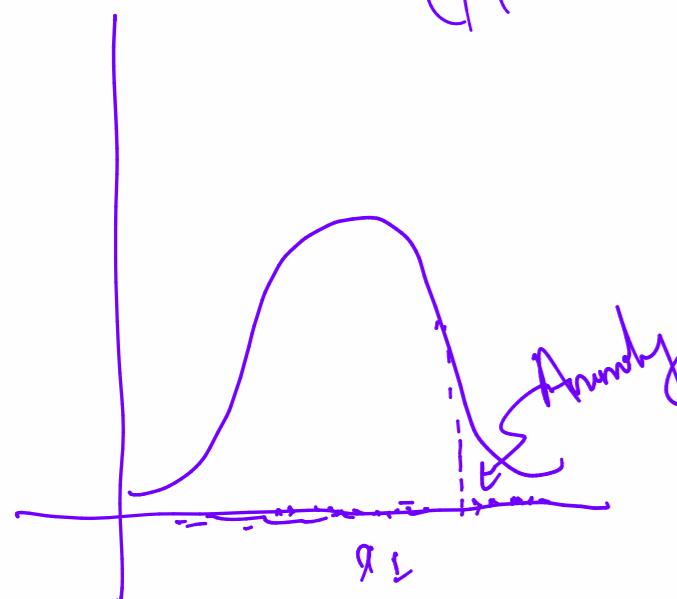
Use
plt.hist(s)
↓
Normalized
np.log(x)

Error analysis for anomaly detection

Want $p(x) \geq \epsilon$ Large for normal examples x .
 $p(x) < \epsilon$ Small for anomalous examples x .

Most common problem

$p(x)$ is comparable for normal and anomalous examples.
($p(x)$ is large for both)



Monitoring computers in a data center

Choose features that might take on unusually large or small values in the event of an anomaly.

Computer Behaving Anomaly?

x_1 = memory use of computer

x_2 = no. of disk access ps/sec

high x_3 = CPU load
low x_4 = network traffic

← Not unusual.

$$x_5 = \frac{\text{CPU load}}{\text{network traffic}}$$

$$x_6 = \frac{C(\text{CPU load})^2}{\text{network traffic}}$$

Deciding feature choice base in $p(x)$.

Making Recommendation

Predicting movie ratings

User rates movies using one to five stars

Movie	Alice(1)	Bob(2)	Carol(3)	Dave(4)
Love at last	5	5	0	0
Romance forever	5	?	0	0
Cute puppies of love	?	4	0	?
Nonstop car chases	0	0	5	4
Swords vs. karate	0	0	5	2

$y_{(i,j)}$ = rating given by user j to movie i
 $\gamma_{(i,j)} = 1$ if user j has rated movie i

$$n_u = 4$$

$$n_m = 5$$

Ratings				
★				
★	★			
★	★	★		
★	★	★	★	
★	★	★	★	★

$$\rightarrow n_u = \text{no. of users}$$

$$\rightarrow n_m = \text{no. of movies}$$

$$\gamma_{(i,j)} = 1 \text{ if user } j \text{ has rated movie } i$$

$$\gamma_{(1,1)} = 1$$

$$\gamma_{(3,2)} = 0$$

Goal: predict
missing ratings for
common items.

Noobdim:

n_u, n_m : num of
users until items.

$r(i,j)$: if user
 j rated movie i , else 0

$y(i,j)$: User j 's rating for movie i .

Predicting movie ratings

User rates movies using one to five stars
zero

Movie	Alice(1)	Bob(2)	Carol(3)	Dave(4)
Love at last	5	5	0	0
Romance forever	5	?	?	0
Cute puppies of love	?	4	0	?
Nonstop car chases	0	0	5	4
Swords vs. karate	0	0	5	?

$$n_u = 4$$

$$r(1,1) = 1$$

$$n_m = 5$$

$$r(3,1) = 0$$

$$y^{(3,2)} = 4$$

$y^{(i,j)}$ = rating given by
user j to movie i
(defined only if $r(i,j) = 1$)

Ratings				
★				
★	★			
★	★	★		
★	★	★	★	
★	★	★	★	★



n_u = no. of users

n_m = no. of movies

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

Start with known features (eg. genres)
Later explore feature-free Methods.