

ANOMALY DETECTION

Lecture 89: Problem Motivation

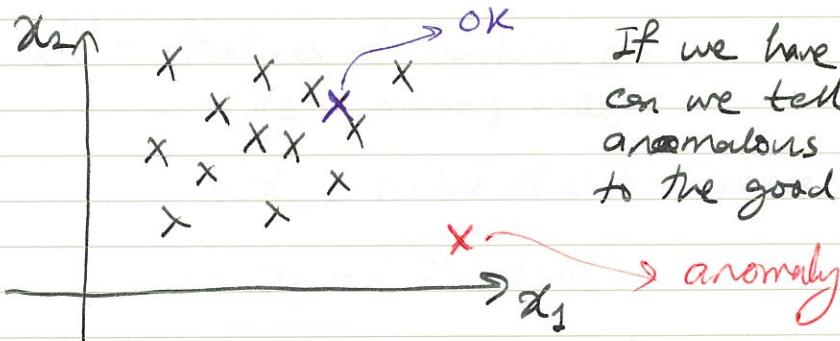
Consider an aircraft manufacturer that wants to make sure no defective aircrafts are given to customers. We can describe an aircraft by a given set of features:

x_1 = heat generated

x_2 = vibration intensity

:

Suppose we plot $x_2 - x_1$ for m aircrafts that we know are good:



If we have a new aircraft, can we tell if it's anomalous by comparing it to the good ones?

More generally, we're given a data set $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$ which we assume is normal (not anomalous), and try to determine if a new example x_{test} is anomalous. We do this by modeling the probability $P(x)$ of an example x occurring. Having that at our disposal, we can make decisions for new examples:

$$\begin{cases} P(x_{\text{test}}) < \varepsilon \Rightarrow \text{Anomalous} \\ P(x_{\text{test}}) \geq \varepsilon \Rightarrow \text{OK} \end{cases}$$

Small threshold

Anomaly detection examples

* Fraud detection:

- $\chi^{(i)}$ = features of user i 's activities:

→ χ_1 = How often does user log in?

→ χ_2 = Number of webpages visited.

→ χ_3 = Number of transactions

→ χ_4 = Number of posts on a forum

→ χ_5 = Typing speed in characters per second

⋮

- Model $P(\chi)$ from data and use it to identify unusual users which have $P(\chi) < \varepsilon$.

* Manufacturing (car aircraft example, for instance)

* Monitoring computers in a data centre:

- $\chi^{(i)}$ = features of machine i :

→ χ_1 = memory usage

→ χ_2 = number of disk accesses / sec

→ χ_3 = CPU load

→ χ_4 = CPU load / network traffic

⋮

- If $P(\chi) < \varepsilon$ for a given machine, we can flag it for inspection since it might be going down.

Lecture 90: Gaussian Distribution

$$P(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

↓

& we write $x \sim N(\mu, \sigma^2)$
 distributed as mean ↓ variance ↓

This is also written as $P(x; \mu, \sigma^2)$ sometimes.

Parameter estimation

Suppose we have a data set $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$ where $x^{(i)} \in \mathbb{R}$ and we suspect $x^{(i)}$ are samples of a normal distribution $N(\mu, \sigma^2)$.

How do we estimate μ & σ^2 ?

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

$$\hat{\sigma}^2 = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \hat{\mu})^2$$

Let's remind ourselves how these estimates are obtained. If $x^{(1)}, \dots, x^{(m)}$ are drawn independently from $N(\mu, \sigma^2)$, then the probability density of their observation is given by:

$$P = \prod_{i=1}^m \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{(x^{(i)} - \mu)^2}{2\sigma^2}\right] = \frac{1}{(2\pi)^{m/2}\sigma^m} e^{-\frac{\sum_{i=1}^m (x^{(i)} - \mu)^2}{2\sigma^2}}$$

Let's pick the values of μ & σ which maximize P .

$$\frac{\partial P}{\partial \mu} = 0 \Rightarrow P \cdot \frac{1}{2\sigma^2} \cdot \sum_{i=1}^m 2(x^{(i)} - \mu) \cdot x - 1 = 0 \Rightarrow \hat{\mu} = \frac{1}{m} \sum_{i=1}^m x^{(i)}$$

Now suppose we pick m samples & compute $\hat{\mu}$, then pick another

m samples & compute $\hat{\mu}$, and so on. Obviously, each time we'd obtain a different number. What would be the mean of these numbers?

Remember that we're assuming $x^{(i)} \sim N(\mu, \sigma^2)$, so:

$$\langle \hat{\mu} \rangle = \frac{1}{m} \sum_{i=1}^m \langle x^{(i)} \rangle = \frac{1}{m} \sum_{i=1}^m \mu = \mu.$$

So $\langle \hat{\mu} \rangle = \mu$, the true mean. For this reason, $\hat{\mu}$ is called an unbiased estimator.

Let's go through the same procedure for σ :

$$\frac{\partial P}{\partial \sigma} = 0 \Rightarrow \frac{1}{(2\pi)^{m/2}} \sigma^{-m-1} \exp \left[-\frac{1}{2\sigma^2} \sum_{i=1}^m (x^{(i)} - \mu)^2 \right]$$

$$+ P \cdot \frac{1}{\sigma^3} \sum_{i=1}^m (x^{(i)} - \mu)^2 = 0$$

$$\Rightarrow \left[-\frac{m}{\sigma} + \frac{1}{\sigma^3} \sum_{i=1}^m (x^{(i)} - \mu)^2 \right] P = 0$$

$$\Rightarrow \hat{\sigma}^2 = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu)^2.$$

Of course, we only have an estimate of μ (i.e. $\hat{\mu}$), so let's take our variance estimator to be: $\hat{\sigma}^2 = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \hat{\mu})^2$.

Is $\hat{\sigma}^2$ unbiased?

$$\hat{\sigma}^2 = \frac{1}{m} \sum_{i=1}^m \left[x^{(i)} - \frac{1}{m} \sum_{i'=1}^m x^{(i')} \right]^2$$

$$= \frac{1}{m} \sum_{i=1}^m [x^{(i)} x^{(i)}] - \frac{2}{m} x^{(i)} \sum_{i'=1}^m x^{(i')} + \frac{1}{m^2} \sum_{i,i''=1}^m x^{(i)} x^{(i'')}$$

$$\hat{\sigma}^2 = \frac{1}{m} \sum_{i=1}^m (\bar{x}^{(i)})^2 - \frac{1}{m^2} \sum_{i,i'=1}^m \bar{x}^{(i)} \bar{x}^{(i')}$$

$$\langle \hat{\sigma}^2 \rangle = \frac{1}{m} \sum_{i=1}^m \langle (\bar{x}^{(i)})^2 \rangle - \frac{1}{m^2} \sum_{i,i'=1}^m \langle \bar{x}^{(i)} \bar{x}^{(i')} \rangle$$

$$\begin{aligned} \langle (\bar{x}^{(i)})^2 \rangle &= \langle (\bar{x}^{(i)} - \mu + \mu)^2 \rangle = \langle (\bar{x}^{(i)} - \mu)^2 + 2\mu(\bar{x}^{(i)} - \mu) + \mu^2 \rangle \\ &= \langle (\bar{x}^{(i)} - \mu)^2 \rangle + 2\mu \langle \bar{x}^{(i)} - \mu \rangle + \mu^2 = \sigma^2 + \mu^2 \end{aligned}$$

$$\langle \bar{x}^{(i)} \bar{x}^{(i')} \rangle = \langle (\bar{x}^{(i)} - \mu + \mu)(\bar{x}^{(i')} - \mu + \mu) \rangle = \sigma^2 S_{ii'} + \mu^2 \quad \left(\begin{array}{l} \text{we've used} \\ \langle (\bar{x}^{(i)} - \mu)(\bar{x}^{(i')} - \mu) \rangle \\ = S_{ii'} \sigma^2 \end{array} \right)$$

$$\begin{aligned} \langle \hat{\sigma}^2 \rangle &= \frac{1}{m} \sum_{i=1}^m (\sigma^2 + \mu^2) - \frac{1}{m^2} \sum_{i,i'=1}^m \sigma^2 S_{ii'} + \mu^2 \\ &= \sigma^2 + \mu^2 - \frac{1}{m^2} (\sigma^2 \cancel{+ m} + \mu^2 \cancel{+ m^2}) \\ &= \sigma^2 \left(1 - \frac{1}{m}\right) = \frac{m-1}{m} \sigma^2 \end{aligned}$$

different examples are independent

We see that $\hat{\sigma}^2$ is not unbiased. That's why sometimes people use $\hat{\sigma}^2 = \frac{1}{m-1} \sum_{i=1}^m (\bar{x}^{(i)} - \hat{\mu})^2$ which would be unbiased. For large values of m this won't make any difference in practice. Andrew Ng claims that in machine learning the biased estimator is used,

i.e. $\hat{\sigma}^2 = \frac{1}{m} \sum_{i=1}^m (\bar{x}^{(i)} - \hat{\mu})^2$.

Lecture 91: Algorithm

Consider the training set $\{x^{(1)}, \dots, x^{(m)}\}$ where $x^{(i)} \in \mathbb{R}^n$. We will assume that each feature is distributed normally & that the different features are not correlated:

$$x_1 \sim N(\mu_1, \sigma_1^2), \dots, x_n \sim N(\mu_n, \sigma_n^2)$$

We can then write $p(x)$ as:

$$P(x) = P(x_1; \mu_1, \sigma_1^2) \times P(x_2; \mu_2, \sigma_2^2) \times \dots \times P(x_n; \mu_n, \sigma_n^2)$$

$$\text{where } P(x_j; \mu_j, \sigma_j^2) = \frac{1}{\sqrt{2\pi}\sigma_j} \exp\left[-\frac{1}{2} \frac{(x_j - \mu_j)^2}{\sigma_j^2}\right]$$

μ_j, σ_j^2 are estimated as follows:

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

Later we will talk about generalizing this result so that correlations are captured, via the multivariate Gaussian. Let's summarize our model of the probability distribution $P(x)$:

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

$$\text{where } P(x_j; \mu_j, \sigma_j^2) = \frac{1}{\sqrt{2\pi}\sigma_j} e^{-\frac{(x_j - \mu_j)^2}{2\sigma_j^2}}$$

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

Then given a new example x_{test} , we say x_{test} is anomalous if $P(x_{\text{test}}) < \varepsilon$ for small ε .

This algorithm makes a few questionable assumptions: we're basically fitting for n uncorrelated Gaussians. What if our features are not normally distributed? What if they're correlated? We'll talk about those issues & more next.

Lecture 92: Developing and Evaluating an Anomaly Detection System

In lecture 67 (see pages 27-28 of Notebook #2) we talked about the importance of having a single real number to evaluate a learning algorithm.

For instance, when picking one feature over another, or deciding what threshold to use in logistic regression for predictions, it'd be useful to have one metric that can tell us whether our learning algorithm is doing better or worse.

In our discussion of supervised learning, we developed some metrics:

cross validation error $J_{\text{cv}} = \frac{1}{2m_{\text{cv}}} \sum_{i=1}^{m_{\text{cv}}} [h_{\theta}(x_{\text{cv}}^{(i)}) - y_{\text{cv}}^{(i)}]^2$ for linear regression, F_1 score (see Lecture 69 - pgs 30-33 of Notebook #2) for classification problems (especially skewed datasets), etc.

How can we evaluate our anomaly detection algorithm in a similar way?

Let's assume our data set is labelled so $y=0$ if normal, $y=1$ if anomalous.

For instance, in our aircraft manufacturing example, we might have some examples of aircrafts with defective engines which had to be fixed. In this case, we can split our data as follows:

* Training set: $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$ where all examples are normal. This makes sense since ^{we} want the probability distribution of non-anomalous examples. Andrew Ng mentions it's okay if a bunch of anomalous examples find their way in the training set, but we should try to keep them out.

* Cross-validation set: $\{(x_{cv}^{(1)}, y_{cv}^{(1)}), \dots, (x_{cv}^{(m)}, y_{cv}^{(m)})\}$.

where we have both anomalous & normal examples.

* Test set: $\{(x_{test}^{(1)}, y_{test}^{(1)}), \dots, (x_{test}^{(m)}, y_{test}^{(m)})\}$ where we have both anomalous & normal examples.

In our aircraft manufacturing example, suppose we have the following data:

* 10,000 good (normal) engines

* 20 flawed engines (anomalous)

This is a typical anomalous/normal ratio in anomaly detection problems. We can construct the above data sets as follows:

* Training set: 6000 good engines

* Cross-validation set: 2000 good engines ($y=0$), 10 anomalous ($y=1$)

* Test set: 2000 good engines ($y=0$), 10 anomalous ($y=1$)

Equipped with this, we can proceed as follows:

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

(ii) On a cross-validation/test set example x , predict:

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

(iii) Use an evaluation metric to determine how well anomalies are detected:

- * True positive, false positive, false negative, true negative

- * Precision / Recall

- * F_1 -score (single #)

Note that the dataset is very skewed, so we should be using precision/recall or F_1 -score.

We can also choose ϵ using the validation set: pick ϵ that results in the highest F_1 -score. What happens if we take ϵ to be very small?

Precision will be high since we're predicting very few anomalies that are likely to be actual anomalies, but recall will be low because of all actual anomalies, we've caught only a few. The opposite is true if ϵ is too large. We can use F_1 -score to pick the optimal ϵ .

Infer want a mix number of positive and negative examples. In the case of spam classification, there are many types of spam emails, but

Lecture 93: Anomaly Detection vs. Supervised Learning

In the previous lecture we assumed we have a labeled dataset to be able to evaluate anomaly detection. If we have a labeled dataset, however, why wouldn't we just use a supervised classification algorithm?

The crucial determinant is the number of positive examples. If we have a very small number of positive example (e.g. 0 - 20), it will be hard for a supervised learning algorithm to learn what positive examples are generally like. This is especially so when there are many different types of anomalies, as is the case with aircraft engine failure. There are so many potential causes of engine failure, that it will be hard for any algorithm to learn from a handful of positive examples what the anomalies look like; future anomalies may look nothing like any of the anomalous examples we've seen so far. Therefore, in such cases, it makes more sense to model the negative examples, of which there are often many, and use that to detect anomalies. As we discussed in the previous lecture, the handful of positive examples can be used to evaluate the anomaly detection algorithm. For supervised learning algorithms to work well, we often want a large number of positive and negative examples. In the case of Spam classification, there are many types of spam emails, but

our training sets often have many spam examples, so it's quite likely that an unseen spam email is similar to some examples in the training set.

Typical examples of anomaly detection vs. supervised learning:

Anomaly detection

- * Fraud detection
- * Manufacturing (e.g. aircraft engines)
- * Monitoring machines in a data center.

Supervised learning

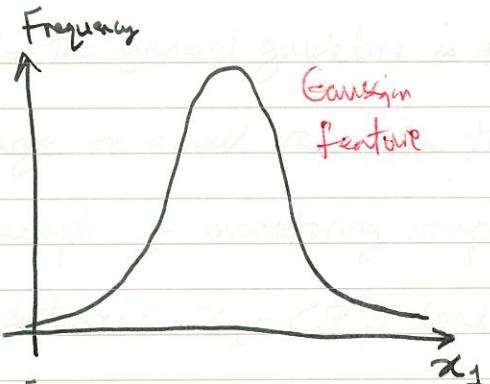
- * Email spam classification
- * Weather prediction (sunny/rainy/etc)
- * Cancer classification

When there's a lot of positive examples, some of these anomaly detection examples could be studied using supervised learning. For instance, if a giant online retailer has many cases of fraudulent behaviour, they could use supervised learning to detect anomalies.

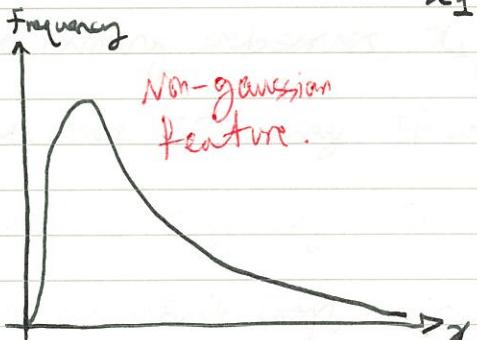
Lecture 94: Choosing What Features to Use

In practice, the performance of the anomaly detection algorithm is affected quite a lot by the choice of features. We modeled our features as independent normally-distributed variables. The Gaussian assumption is a bold one and may not be true. Therefore, it's a good idea to draw the histogram of values of a given feature across our training set:





If we get a plot like this, then we know our assumption is probably okay.



If we get a plot like this, we should transform our feature so that it becomes more Gaussian.

For instance; we might try: $x_1 \rightarrow \log x_1$,

which would do the trick if x_1 has a log-normal distribution.

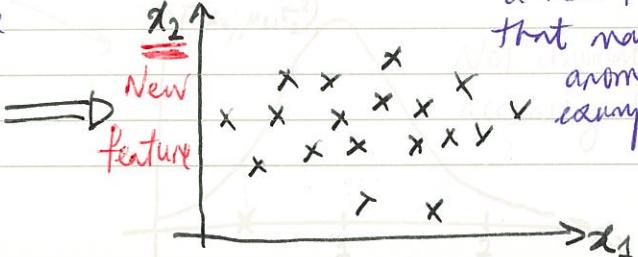
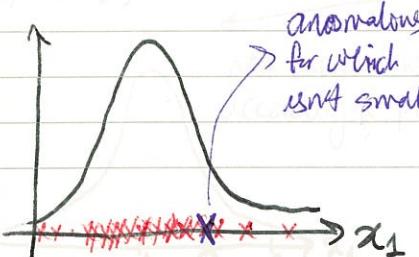
This is, of course, just one example. Here are some others:

$$x_2 \rightarrow \sqrt{x_2}, \quad x_3 \rightarrow x_3^3, \quad x_4 \rightarrow \log(x_4 + c), \text{ etc.}$$

For anomaly detection, we want our modeled probability distribution to have the following behaviour:

$$\begin{cases} P(x) \text{ large for normal examples } x. \\ P(x) \text{ small for anomalous examples } x. \end{cases}$$

A common problem, however, is that $P(x)$ might end up being comparable for normal and anomalous examples. In such cases, we should study the anomalous examples for which $P(x)$ is not small and try to come up with features that make them stand out,

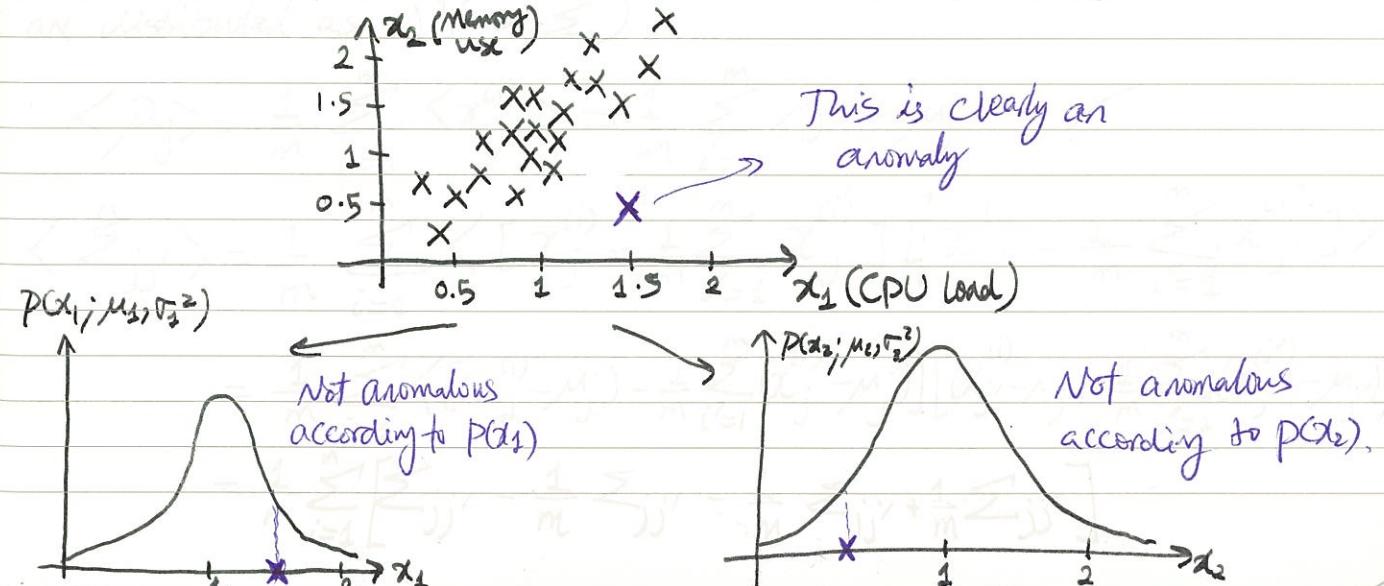


$\times \rightarrow$ We've created a new feature x_2 that makes our anomalous example stand out.

So the general guideline is to choose features that might take on unusually large or small values in the event of an anomaly. Let's go back to the example of monitoring computers in a data center. Consider the two features: $x_1 = \text{CPU load}$ & $x_2 = \text{Network traffic}$. If the computers are running web servers, x_1 & x_2 are correlated: the more network traffic, the more CPU usage. If we want to distinguish computers for which CPU load is not growing with network traffic, for instance because they're stuck in an infinite loop, we can create a new feature $x_3 = x_1 / x_2$, which takes unusually high values if CPU usage is high despite modest or low network traffic.

Lecture 95: Multivariate Gaussian Distribution

The anomaly detection algorithm we described in Lecture 91 assumes different features are not correlated. This can cause issues:



Because our model is assuming x_1 & x_2 are not correlated, it will not detect the anomaly. To get around this, we'll model $x \in \mathbb{R}^n$ as a multivariate Gaussian parametrised by $\mu \in \mathbb{R}^n$ (mean) & $\Sigma \in \mathbb{R}^{n \times n}$ (covariance matrix):

$$-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu)$$

$$P(x; \mu, \Sigma) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} e^{-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu)}$$

determinant of Σ .

It can be shown that $\begin{cases} \sum_{jj'} = \langle (x_j - \mu_j)(x_{j'} - \mu_{j'}) \rangle \\ \mu_j = \langle x_j \rangle \end{cases}$

Let's go through our sample estimate arguments (see Lecture 9D) again:

$$\hat{\mu}_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)}$$

$$\hat{\Sigma}_{jj'} = \frac{1}{m} \sum_{i=1}^m (x_j^{(i)} - \hat{\mu}_j)(x_{j'}^{(i)} - \hat{\mu}_{j'})$$

Let's check if $\hat{\mu}$ & $\hat{\Sigma}$ are unbiased estimates, remembering that $x^{(i)}$ are distributed as $\mathcal{N}(\mu, \Sigma)$.

$$\langle \hat{\mu}_j \rangle = \frac{1}{m} \sum_{i=1}^m \langle x_j^{(i)} \rangle = \frac{1}{m} \sum_{i=1}^m \mu_j = \mu_j \quad \checkmark$$

$$\begin{aligned} \langle \hat{\Sigma}_{jj'} \rangle &= \frac{1}{m} \sum_{i=1}^m \left\langle \left[x_j^{(i)} - \frac{1}{m} \sum_{i'=1}^m x_j^{(i')} \right] \left[x_{j'}^{(i)} - \frac{1}{m} \sum_{i''=1}^m x_{j''}^{(i'')} \right] \right\rangle \\ &= \frac{1}{m} \sum_{i=1}^m \left\langle \left[(x_j^{(i)} - \mu_j) - \frac{1}{m} \sum_{i'=1}^m (x_j^{(i')} - \mu_j) \right] \left[(x_{j'}^{(i)} - \mu_{j'}) - \frac{1}{m} \sum_{i''=1}^m (x_{j''}^{(i'')} - \mu_{j''}) \right] \right\rangle \\ &= \frac{1}{m} \sum_{i=1}^m \left[\sum_{j,j'} \left(x_j^{(i)} - \mu_j \right) \left(x_{j'}^{(i)} - \mu_{j'} \right) - \frac{1}{m} \sum_{i=1}^m \sum_{j,j'} \left(x_j^{(i)} - \mu_j \right) \left(x_{j'}^{(i)} - \mu_{j'} \right) \right] \end{aligned}$$

where we've used that $\langle (x_j^{(i)} - \mu_j) (x_{j'}^{(i')} - \mu_{j'}) \rangle = \sum_{jj'} \underline{s_{jj'}}$
 different examples
 are independent

$\langle \hat{\Sigma}_{jj'} \rangle = \frac{m-1}{m} \sum_{jj'} \underline{s_{jj'}}$, much as it was the case in the case
 of a single Gaussian. For large values of m , $\frac{m-1}{m} \approx 1$ & this won't
 matter much.

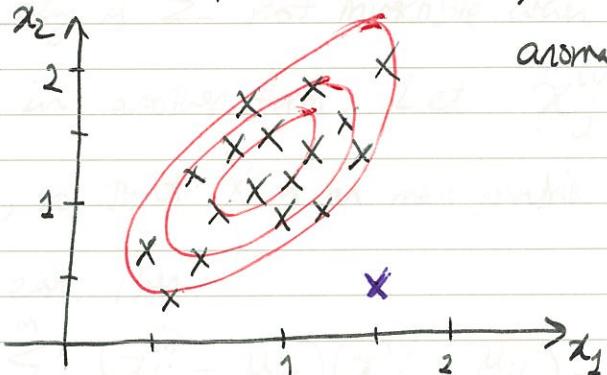
Lecture 96: Anomaly Detection Using the Multivariate Gaussian Distribution

Using the multivariate Gaussian to model $p(x)$, we have:

$$p(x) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp\left[-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu)\right]$$

where $\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)}$ & $\Sigma = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu)^T (x^{(i)} - \mu)^T$

For our example in the previous lecture, $p(x)$ would now be small for the



anomalous example.

The model introduced in lecture 91 is a special case of the model above, where $\Sigma_{jj'} = \sigma_j^2 \delta_{jj'}$. As we've already mentioned, the multivariate

Gaussian model automatically captures correlations, whereas for the original model we need to manually create new features (like $\frac{\text{CPU load}}{\text{memory usage}}$) to handle correlations. The original model, however, is computationally cheaper and scales better for large values of n . The reason is that we need Σ^{-1} , which typically scales as n^3 . For $n \gtrsim 10,000$, inverting Σ becomes too computationally expensive. We faced the same dilemma in linear regression, when choosing between the Normal Equation and Gradient Descent (see Lecture 24, pgs 22-24 of Notebook #1).

Another case where we cannot use the multivariate Gaussian model is when $m < n$, since in that case Σ is not invertible. We will soon prove this, but it shouldn't be that hard to believe: we're fitting about $\frac{n^2}{2}$ parameters, so we need a lot of data.

Okay, so why is Σ not invertible when $m < n$? Let's first rewrite Σ in another form. Let $\hat{x}_j^{(i)} = x_j^{(i)} - \mu_j$ & $\hat{X}_{ij} = \hat{x}_j^{(i)}$, so that \hat{X} is an $m \times n$ matrix where every feature (ie column) has zero mean. Then:

$$\begin{aligned}\Sigma_{jj'} &= \frac{1}{m} \sum_{i=1}^m (x_j^{(i)} - \mu_j)(x_{j'}^{(i)} - \mu_{j'}) = \frac{1}{m} \sum_{i=1}^m \hat{x}_j^{(i)} \hat{x}_{j'}^{(i)} \\ &= \frac{1}{m} \sum_{i=1}^m \hat{X}_{ij} \hat{X}_{ij'} = \frac{1}{m} [\hat{X}^T \hat{X}]_{jj'} \Rightarrow \Sigma = \hat{X}^T \hat{X}.\end{aligned}$$

This is actually what we have already seen in PCA (See page 35 for instance). Let $A = \hat{X}/\sqrt{m}$ so that $\Sigma = A^T A$.

Of course, Σ is symmetric & positive semi-definite (see page 32 for the proof) and as a result can be diagonalized and all of its eigenvalues are greater or equal to zero.

$$\Sigma v^{(i)} = \lambda_i v^{(i)} \quad i=1, \dots, n$$

$$\text{where } v^{(i)} \cdot v^{(j)} = \delta_{ij} \text{ & } \lambda_i \geq 0.$$

Now A is an $m \times n$ matrix, so $W \equiv A A^T$ is an $m \times m$ symmetric & positive definite matrix.

$$W w^{(i)} = \xi_i w^{(i)} \quad i=1, \dots, m$$

$$\text{where } w^{(i)} \cdot w^{(j)} = \delta_{ij} \text{ & } \xi_i \geq 0.$$

Claim I: $A v^{(i)}$ is an eigenvector of W with eigenvalue λ_i if $\lambda_i > 0$.

$$\begin{aligned} \text{Proof. } W(A v^{(i)}) &= A A^T (A v^{(i)}) = A (\bar{A}^T A) v^{(i)} = A \Sigma v^{(i)} \\ &= A (\lambda_i v^{(i)}) = \lambda_i (A v^{(i)}). \end{aligned}$$

Claim II: $A^T w^{(i)}$ is an eigenvector of Σ with eigenvalue ξ_i if $\xi_i > 0$.

$$\begin{aligned} \text{Proof. } \Sigma (A^T w^{(i)}) &= A^T A (A^T w^{(i)}) = A^T (A A^T) w^{(i)} = A^T W w^{(i)} \\ &= A^T (\xi_i w^{(i)}) = \xi_i (A^T w^{(i)}). \end{aligned}$$

$m = \min(n, m)$ non-zero eigenvalues.

What about $A\gamma^{(i)}$ when $\lambda_i = 0$?

$$(A\gamma^{(i)})^T (A\gamma^{(i)}) = \gamma^{(i)T} A^T A \gamma^{(i)} = \gamma^{(i)T} \Sigma V^{(i)} = \lambda_i = 0$$

$$\Rightarrow A\gamma^{(i)} = 0.$$

Similarly with $A^T w^{(i)}$ when $\xi_i = 0$:

$$(A^T w^{(i)})^T (A^T w^{(i)}) = w^{(i)T} A A^T w^{(i)} = w^{(i)T} \Sigma W w^{(i)} = \xi_i = 0$$

$$\Rightarrow A^T w^{(i)} = 0.$$

Claim III: The maximum number of non-zero eigenvalues of both Σ & W

is $\min(m, n)$. Also, Σ & W share the same set of eigenvalues.

Proof. Assume $m > n$ & $\xi_i > 0 \forall i=1, \dots, m$. Let $\Xi^{(i)} = \frac{1}{\sqrt{\xi_i}} A^T w^{(i)}$.

$$\begin{aligned} \text{Then } \Xi^{(i)T} \Xi^{(j)} &= \frac{1}{\sqrt{\xi_i} \sqrt{\xi_j}} (A^T w^{(i)})^T (A^T w^{(j)}) = \frac{1}{\sqrt{\xi_i \xi_j}} w^{(i)T} (A A^T) w^{(j)} \\ &= \frac{1}{\sqrt{\xi_i \xi_j}} w^{(i)T} W w^{(j)} = \frac{\xi_j}{\sqrt{\xi_i \xi_j}} g_{ij} = g_{ij} \quad \forall i, j = 1, \dots, m. \end{aligned}$$

We've shown that $\Xi^{(1)}, \dots, \Xi^{(m)}$ form an orthonormal basis. Remember that $\Xi^{(i)} \in \mathbb{R}^n$, and since $m > n$, this is impossible! Therefore W can only have n non-zero eigenvalues at maximum. In this case, $n = \min(n, m)$.

Assume $n > m$ & $\lambda_i > 0 \forall i=1, \dots, n$. Let $\Xi^{(i)} = \frac{1}{\sqrt{\lambda_i}} A \gamma^{(i)}$.

$$\text{Then } \Xi^{(i)T} \Xi^{(j)} = \frac{1}{\sqrt{\lambda_i \lambda_j}} \gamma^{(i)T} \underbrace{(A^T A)}_{\Sigma} \gamma^{(j)} = g_{ij}. \text{ Again, we cannot have}$$

n orthonormal vectors in \mathbb{R}^m when $n > m$, therefore Σ can at most have $m = \min(n, m)$ non-zero eigenvalues.

That Σ & W share the same eigenvalues follows from claims I & II.

If λ_i is a non-zero eigenvalue of Σ , it's also an eigenvalue of W . But can W have more non-zero eigenvalues than Σ ? No, because every non-zero eigenvalue of W is also a non-zero eigenvalue of Σ .

So, to get back to why Σ is non-invertible when $m < n$: this is because it will definitely have a zero eigenvalue, since it can have at most m non-zero eigenvalues.

With the results we've proven so far, it would be a shame not to discuss the Singular-Value Decomposition (SVD).

Claim IV: Every $m \times n$ matrix A can be written as:

$$A = UDV^T \text{ where: } \begin{aligned} &+ U \text{ is } m \times m \text{ & satisfies } U^T U = U U^T = I_m \\ &+ V \text{ is } n \times n \text{ & satisfies } V^T V = V V^T = I_n \\ &+ D \text{ is } m \times n \text{ & satisfies } D_{ij} = \alpha_i S_{ij} \end{aligned}$$

$$\text{where } i=1, \dots, m, j=1, \dots, n \text{ & } \alpha_i \geq 0.$$

Proof: Let $\lambda_1, \dots, \lambda_k$ ($k \leq \min(m, n)$) denote the strictly positive eigenvalues of $\Sigma \succeq W$. Let $v^{(1)}, \dots, v^{(k)}$ be the corresponding eigenvectors of Σ .

Define: $V = \left[\begin{bmatrix} v^{(1)} \\ \vdots \\ v^{(k)} \end{bmatrix} \dots \begin{bmatrix} v^{(k)} \\ \vdots \\ v^{(n)} \end{bmatrix} \right] = \left[\begin{bmatrix} v^{(1)} \\ \vdots \\ v^{(k+1)} \\ \vdots \\ v^{(n)} \end{bmatrix} \dots \begin{bmatrix} v^{(k)} \\ \vdots \\ v^{(n)} \end{bmatrix} \right]$

$$\sum \lambda_i v^{(i)} = \lambda_i v^{(i)}$$

$$\Leftarrow v^{(i)} \cdot v^{(j)} = S_{ij}$$

$$\Leftarrow \lambda_1, \dots, \lambda_k > 0$$

$$\Leftarrow \lambda_{k+1} = \dots = \lambda_n = 0$$

*orthonormal eigenvectors
corresponding to zero eigenvalues of Σ .*

Similarly, let $w^{(1)} = \frac{1}{\sqrt{\lambda_1}} A y^{(1)}$, ..., $w^{(k)} = \frac{1}{\sqrt{\lambda_k}} A y^{(k)}$

& pick $w^{(k+1)}, \dots, w^{(m)}$ to be orthonormal eigenvectors of W corresponding to zero eigenvalue ξ . $w^{(i)} \cdot w^{(j)} = \delta_{ij}$ & $W W^{(i)} = \xi_i w^{(i)}$

where as we've shown $\xi_1 = \lambda_1$, $\xi_2 = \lambda_2$, ..., $\xi_k = \lambda_k$ & $\xi_{k+1} = \dots = \xi_m = 0$.

Define: $V = \begin{bmatrix} [w^{(1)}] & \dots & [w^{(k)}] & [w^{(k+1)}] & \dots & [w^{(m)}] \end{bmatrix}$

$$= \frac{1}{\sqrt{\lambda_1}} A y^{(1)} \quad = \frac{1}{\sqrt{\lambda_k}} A y^{(k)}$$

eigenvectors corresponding to zero eigenvalue of W .

of course, V^T is nanc & V is marn.

Let $D_{ij} = \sqrt{\lambda_i} \delta_{ij}$, $i=1, \dots, m$ & $j=1, \dots, n$.

So $D_{11} = \sqrt{\lambda_1}, \dots, D_{kk} = \sqrt{\lambda_k}$ & all other entries of D are zeros.

Then:

$$V D V^T = \begin{bmatrix} \left[\frac{A y^{(1)}}{\sqrt{\lambda_1}} \right] & \dots & \left[\frac{A y^{(k)}}{\sqrt{\lambda_k}} \right] & [w^{(k+1)}] & \dots & [w^{(m)}] \end{bmatrix}_{m \times n} \begin{bmatrix} \sqrt{\lambda_1} & & & & & \\ & \sqrt{\lambda_2} & & & & \\ & & \ddots & & & \\ & & & \sqrt{\lambda_k} & & \\ & & & & \ddots & \\ & & & & & \sqrt{\lambda_m} \end{bmatrix}_{m \times m} \begin{bmatrix} [V^{(1)T}]^T & & & & & \\ & \vdots & & & & \\ & & [V^{(k)T}]^T & & & \\ & & & [V^{(k+1)T}]^T & & \\ & & & & \ddots & \\ & & & & & [V^{(m)T}]^T \end{bmatrix}_{n \times n}$$

$$= \begin{bmatrix} \left[\frac{A y^{(1)}}{\sqrt{\lambda_1}} \right] & \dots & \left[\frac{A y^{(k)}}{\sqrt{\lambda_k}} \right] & [w^{(k+1)}] & \dots & [w^{(m)}] \end{bmatrix}_{m \times n} \begin{bmatrix} [\sqrt{\lambda_1} V^{(1)T}]^T & & & & & \\ & \vdots & & & & \\ & & [\sqrt{\lambda_k} V^{(k)T}]^T & & & \\ & & & [\sqrt{\lambda_{k+1}} V^{(k+1)T}]^T & & \\ & & & & \ddots & \\ & & & & & [\sqrt{\lambda_m} V^{(m)T}]^T \end{bmatrix}_{m \times n}$$

Carrying out the matrix multiplication leads to :

$$UDV^T = A\bar{V}\bar{V}^T \text{ where } \bar{V} = \left[\begin{bmatrix} v^{(1)} \\ \vdots \\ v^{(k)} \end{bmatrix} \dots \begin{bmatrix} v^{(k+1)} \\ \vdots \\ v^{(n)} \end{bmatrix} \right]_{n \times k}$$

Let $\tilde{V} = \left[\begin{bmatrix} v^{(k+1)} \\ \vdots \\ v^{(n)} \end{bmatrix} \dots \begin{bmatrix} v^{(n)} \end{bmatrix} \right]_{n \times (n-k)}$ be the collection of eigenvectors with zero eigenvalues.

It then follows from $VV^T = I_n$ that $\bar{V}\bar{V}^T + \tilde{V}\tilde{V}^T = I_n$.

Multiplying by A from left: $A\bar{V}\bar{V}^T + A\tilde{V}\tilde{V}^T = A$. But $A\tilde{V} = 0$

since $A v^{(k+1)} = A v^{(k+2)} = \dots = A v^{(n)} = 0$ (see page 64), so: $A\bar{V}\bar{V}^T = A$.

This concludes the proof: $UDV^T = A$. Note that we also showed that σ_i , which are called the singular values of A , are given by positive eigenvalues of $A^T A$ & $A A^T$.

This was a long detour! Back to anomaly detection. Andrew Ng mentions that he would use the multivariate Gaussian model if $m \geq 10n$, just to be safe. Redundant feature could also cause Σ to be non-invertible.