

WEEK - 9

⇒ Bad use of PCA → To prevent overfitting



⇒ ANOMALY DETECTION (Unsupervised)

→ Fraud Detection

- $x^{(i)}$ = features of user i 's activities
- Model $p(x)$ from data
- Identify unusual users by checking which have $p(x) < \epsilon$

→ Manufacturing

- Density estimation → We make a cluster and if the test example is in the cluster, it's OK. otherwise it's an anomaly.

→ Dataset: $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$
New engine: x_{test}

→ Monitoring computers in a data center

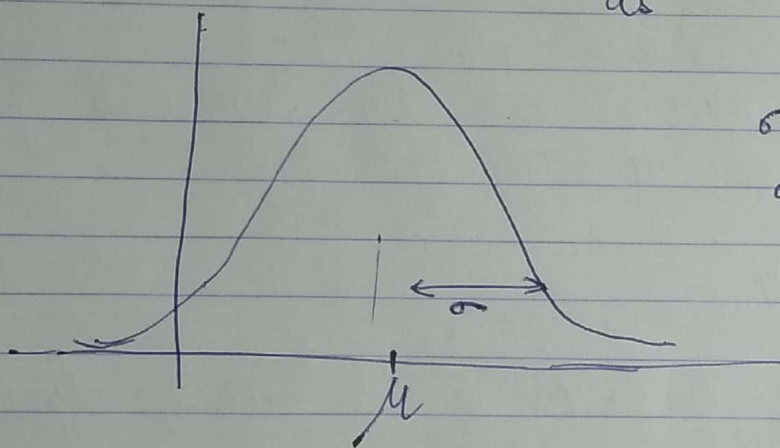
- $x^{(i)}$ = features of machine i
- x_1 = memory use, x_2 = number of disk accesses/sec,
- x_3 = CPU load
- x_4 = CPU load / network traffic divided by.

⇒ Gaussian (Normal) Distribution

Say, $x \in \mathbb{R}$. If x is a distributed Gaussian with mean μ , variance (σ^2)

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

“distributed as” stands for Normal



σ = standard deviation
 σ^2 = variance

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

$$\Rightarrow \begin{aligned} \mu &= ((1/m) * \text{sum}(X))^1; \\ \mu &= \frac{1}{m} \sum_{i=1}^m x^{(i)} \end{aligned} \quad , \quad \begin{aligned} \sigma^2 &= \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu)^2 \end{aligned}$$

$$\Rightarrow \sum_{i=1}^n i = 1 + 2 + 3 + \dots + n \quad (\text{Add})$$

$$\prod_{j=1}^n j = 1 \times 2 \times 3 \times \dots \times n \quad (\text{multiply})$$

$$x_1 \sim \mathcal{N}(\mu_1, \sigma_1^2)$$

⇒ Density of estimation —

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

→ Each example is $x \in \mathbb{R}^n$.

⇒ $P(x) =$ ~~$P(x_1, \mu_1, \sigma_1^2)$~~ Probability of features
 $= P(x_1, \mu_1, \sigma_1^2) * P(x_2, \mu_2, \sigma_2^2) \dots P(x_n, \mu_n, \sigma_n^2)$

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

⇒ Anomaly Detection Algorithm

1) = Choose features x_i that you think might be indicative anomalous examples.

2) = For examples $\{x^{(1)}, x^{(2)} \dots x^{(m)}\}$.
 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 $p(x) < \epsilon$

⇒ ~~How~~ We can use cross validation set to choose parameter ϵ (epsilon)

⇒ When to use anomaly detection and when to use Supervised Learning

Anomaly Detection

- 1) = Very small no. of positive examples ($y=1$). → for (0-50 no) & large no. of -ve examples ($p(x)$)
- 2) = Many different types of anomalies
- 3) = Future anomalies may look nothing like any of the anomalous examples we have seen.

Supervised Learning

- 1) = Large no. of positive and negative examples.
- 2) = Enough +ve examples for algorithm to get a sense of what +ve examples are like.
- 3) = Future +ve examples likely to be similar to ones in training set.

⇒ Anomaly detection usually works good on skewed datasets.

⇒ It is a good habit to make the data gaussian before feeding it to the algo.

⇒ Multivariate Gaussian Distribution

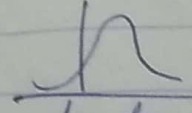
→ $x \in \mathbb{R}^n$. Don't model $p(x_1), p(x_2) \dots$ etc separately.

→ Model $p(x)$ all in one go

Parameters: $\mu \in \mathbb{R}^n$, $\Sigma \in \mathbb{R}^{n \times n}$ (covariance matrix)

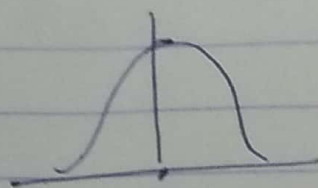
$$p(x; \mu, \Sigma) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)\right)$$

$|\Sigma| \rightarrow$ determinant of Σ / $\det(\text{Sigma})$

→ The area under the bell curve  is always 1 (probability amounts to 1).

⇒ $\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$

$$\Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$



centered at $(0,0) \rightarrow \mu$.

⇒ for varying the value of Σ , the curve flattens or changes its shape

⇒ for $\Sigma = \begin{bmatrix} 0.6 & 0 \\ 0 & 0.6 \end{bmatrix}$ → curves become steeper & taller

⇒ for $\Sigma = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}$ → curves becomes flat & small

⇒ for $\Sigma = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix}$ → curve steepens from N-W & S-E as value increases after 0

⇒ for $\Sigma = \begin{bmatrix} 1 & -0.5 \\ -0.5 & 1 \end{bmatrix}$ → curve steepens from NE & SW as value decreases under 0

⇒ for $\Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 0.6 \end{bmatrix}$ → curve flattens from x_2 side

⇒ Parameter fitting →

Given training set $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$

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

⇒ Anomaly detection with Multivariate Gaussian

1) = Fit model $p(x)$ by setting

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

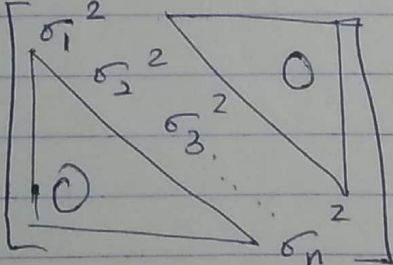
$$\Sigma = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu)(x^{(i)} - \mu)^T$$

2) = Given a new example x , compute

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

Flag an anomaly if $p(x) < \epsilon$

⇒ The original gaussian model is same as the multivariate gaussian ^{not} with a ~~constraint~~ constraint.

where $\Sigma =$  must have 0 everywhere except the diagonal

⇒ The contours of multivariate gaussian are axis aligned in the case of old model → The old model is a special case of multivariate gaussian.

⇒ Original Model

$$p(x_1; \mu_1; \sigma_1^2) \times \dots \times p(x_n; \mu_n; \sigma_n^2)$$

→ Manually create features to capture anomalies where x_1, x_2 take unusual combinations of values

⇒ Computationally cheaper (scales better to large n)

→ OK even if m training set size is small

Multivariate Gaussian

$$p(x; \mu; \Sigma) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp(\dots)$$

→ Automatically captures correlations between features

⇒ Computationally more expensive

⇒ Must have $m \geq n$, or else Σ is non-invertible. ($m \geq 10n$)

⇒ If Σ is non-invertible then —

1) = $m \leq n$

2) = redundant features (similar features)

RECOMMENDER SYSTEMS

→ n_u = no. of users.

→ n_m = no. of movies

→ m_j ⇒ no. of movies rated by user.

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

→ $y(i, j)$ = rating given by user j to movie i .
(depends only if $r(i, j) = 1$)

⇒ To learn $\theta^{(j)}$ (for user j)

$$\min_{\theta^{(j)}} \frac{1}{2m^{(j)}} \sum_{i: r(i, j)=1} \left((\theta^{(j)})^T x^{(i)} - y^{(i, j)} \right)^2$$

(all movies rated by j)

read as → summation of all values of i ,
so the $r(i, j)$ is equal to 1.

$$+ \frac{\lambda}{2m^{(j)}} \sum_{k=1}^n \left(\theta_k^{(j)} \right)^2$$

⇒ To learn $\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(n_u)}$:

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

$$J(\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(n_u)})$$

$$-\frac{\partial}{\partial \theta} (J(\theta^1, \theta^2, \dots, \theta^{n_u}))$$

⇒ gradient descent update

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

for $k=0$

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

(for $k \neq 0$)

⇒ given $\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(n)}$, to learn $x^{(i)}$:

~~min~~

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

sum of all users who
rated the movie i

⇒ given $\theta^{(1)}, \dots, \theta^{(n_u)}$, to learn $x^{(1)}, \dots, x^{(n_m)}$:

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

⇒ Collaborative Filtering Algorithm

→ Given $x^{(1)}, \dots, x^{(n_m)}$ (and movie ratings),
can estimate $\theta^{(1)}, \dots, \theta^{(n_u)}$

→ Given $\theta^{(1)}, \dots, \theta^{(n_u)}$
can estimate $x^{(1)}, \dots, x^{(n_m)}$

→ So what we solve first →

→ we can guess random values of θ
and then solve for x , then taking
those values, we can ~~also~~ get a
better θ , and use this θ for new x
and so on.

And we converge to reasonable features

⇒ Minimizing $x^{(1)}, \dots, x^{(n_m)}$ & $\theta^{(1)}, \dots, \theta^{(n_u)}$
simultaneously.

$$J(x^{(1)}, \dots, x^{(n_m)}, \theta^{(1)}, \dots, \theta^{(n_u)}) = \frac{1}{2} \sum_{(i,j): r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y(i,j))^2$$

$$+ \frac{\lambda}{2} \sum_{i=1}^{n_m} \sum_{k=1}^n (x_k^{(i)})^2 + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^n (\theta_k^{(j)})^2$$

⇒ Basically adding both the equations.

$$\Rightarrow \begin{aligned} X\text{-grad} &= (R * (X^T \Theta - Y)) * \Theta \\ Y\text{-grad} &= (R * (X^T \Theta - Y))^T * X; \end{aligned}$$

8.

⇒ By using the above formula, we do not need to solve for θ, θ' and then for x , again & again.

We can minimize them simultaneously w.r.t each parameter.

⇒ Algo

- 1) = Initialize $x^{(1)} \dots x^{(n_m)}, \theta^{(1)} \dots \theta^{(n_u)}$ to small random values. (this ensures symmetry breaking and ensures the algo learns features $x^{(1)} \dots x^{(n_m)}$ that are different from each other)
- 2) = Minimize $J(x^{(1)}, \dots, x^{(n_m)}, \theta^{(1)}, \dots, \theta^{(n_m)})$ using gradient descent (or an advanced optimization algo). Eg. for every $j=1, \dots, n_u$
 $i=1, \dots, n_m$

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

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

- 3) = For a user with parameters θ and a movie with (learned) features x , predict a star rating of $\theta^T x$.

⇒ Collaborative Filtering

Predicted ratings

$$Y = \begin{bmatrix} 5 & 5 & 0 & 0 \\ 5 & ? & ? & 0 \\ ? & 4 & 0 & ? \\ 0 & 0 & 5 & 4 \\ 0 & 0 & 5 & 0 \end{bmatrix} \quad \begin{bmatrix} (\theta^{(1)})^T x^{(1)} & (\theta^{(2)})^T x^{(1)} & \dots & (\theta^{(n_u)})^T x^{(1)} \\ (\theta^{(1)})^T x^{(2)} & (\theta^{(2)})^T x^{(2)} & \dots & (\theta^{(n_u)})^T x^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ (\theta^{(1)})^T x^{(n_m)} & (\theta^{(2)})^T x^{(n_m)} & \dots & (\theta^{(n_u)})^T x^{(n_m)} \end{bmatrix}$$

$$\Rightarrow \text{Let } X = \begin{bmatrix} -(x^{(1)})^T \\ -(x^{(2)})^T \\ \vdots \\ -(x^{(n_m)})^T \end{bmatrix}, \quad \Theta = \begin{bmatrix} -(\theta^{(1)})^T \\ -(\theta^{(2)})^T \\ \vdots \\ -(\theta^{(n_u)})^T \end{bmatrix}$$

$$\Rightarrow X\Theta^T = \begin{bmatrix} (x^{(1)})^T \theta^{(1)} & \dots & (x^{(1)})^T \theta^{(n_u)} \\ \vdots & \ddots & \vdots \\ \theta(x^{(n_m)})^T \theta^{(1)} & \dots & \theta(x^{(n_m)})^T \theta^{(n_u)} \end{bmatrix}$$

Aka → Low Rank Matrix Factorization.

⇒ To find similar movies, find the movies with smallest $\|x^{(i)} - x^{(j)}\|$, $x^{(i)}$ being your initial movie.

⇒ If a new user is added and it has rated no movies; so instead of giving 0 reviews to every movie, we can mean normalize each row and gives that value as rating.