

Density Estimation

Anomaly Detection

Dataset: $\{x^{(1)}, \dots, x^{(n)}\}$

Is x_{test} anomalous?

Model $p(x)$

$$p(x_{test}) < \varepsilon \rightarrow \text{flag anomaly}$$

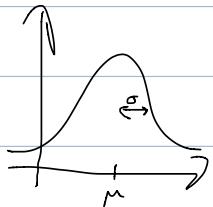
$$p(x_{test}) \geq \varepsilon \rightarrow \text{bk}$$

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

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

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

distribution



$$\begin{aligned} p(x; \mu, \sigma^2) \\ = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \end{aligned}$$

Parameter estimation

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

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

Anomaly Detection Algorithm

+ training set: $\{x^{(1)}, \dots, x^{(n)}\}$

Each example is $x \in \mathbb{R}^n$ $x_i \sim N(\mu_i, \sigma_i^2)$

$p(x)$

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

$$= p(x_1; \mu_1, \sigma_1^2) p(x_2; \mu_2, \sigma_2^2) \cdots p(x_n; \mu_n, \sigma_n^2) \Leftarrow$$

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

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

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

$$\mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)}$$
$$\sigma_j^2 = \frac{1}{m} \sum_{i=1}^m (x_j^{(i)} - \mu_j)^2$$
$$p(x; \mu, \sigma^2) \quad \mu_1, \dots, \mu_n \quad \mu = \begin{bmatrix} \mu_1 \\ \vdots \\ \mu_n \end{bmatrix} = \frac{1}{m} \sum_{i=1}^m x^{(i)}$$

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

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

Anomaly if $p(x) < \epsilon$



Developing and Evaluating an ADS

Importance of real number evaluation

Some labeled data, of anomalous & non-anomalous examples

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

CV & test DO have anomalous examples

$$y = 0 \sim \text{anom}$$

$$y = 1 \sim \text{non-anom}$$

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

On a cross validation/test example x , predict:

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

Possible evaluation metrics:

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

- Precision/Recall

- F₁-Score

(Can also use cross validation set to choose parameter ϵ)

Anomaly Detection vs. Supervised Learning

Very small # of positive examples ($y=1$) ($0 \geq 0$ is common) AD

Large # of negative ($y=0$) examples

Many different "types" of anomalies. H_n , for any algo to learn from (+) examples, what anomalies look like

Large # of positive & negative examples

Enough (+) for algo to get sense of what (+) Ex are like, future will be similar

Choosing What Features to Use

Error analysis for AD

Want $p(x)$ large for normal examples x

$p(x)$ small for anomalous examples x

Most common problem:

$p(x)$ is comparable (say, both 1/2) for normal and anomalous examples

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

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)

$$p(x; \mu, \Sigma)$$

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

Σ is positive definite

A1) Using Multivariate Gaussian Distribution

parameters μ, Σ

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

Parameter fitting:

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

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

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

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

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

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

2. Given a new example x , compute

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

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

Relationship to original model

$$\text{Original model: } p(x) = p(x_1; \mu_1, \sigma_1^2) \times \dots \times p(x_n; \mu_n, \sigma_n^2)$$

Corresponds to multivariate Gaussian

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

$$\text{where } \Sigma = \begin{bmatrix} \sigma_1^2 & & \\ & \sigma_2^2 & \\ & & \ddots & \sigma_n^2 \end{bmatrix}$$

Original Use Case

- used more often

Manually create features to capture anomalies where x_1, x_2 take unusual combinations of values
Computationally cheaper (alternatively, slow to large n)

OK even if m (training set size) is small

Multivariate linear

Automatically captures correlation between features

Computationally more expensive

$$\begin{aligned}x_1 &= \cancel{x}_2 \\ \cancel{x}_3 &= x_4\end{aligned}$$

Must have $m > n$, or else Σ is non-invertible \Leftarrow practice $m \geq 10n$

Recommender Systems

Problem Formulation

$$n_u = \# \text{users}$$

$$n_m = \# \text{movies}$$

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

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

(Content based) Recommendations

$$\Theta^{(j)} \in \mathbb{R}^{n+1}$$

For each user j , know a parameter, $\Theta^{(j)} \in \mathbb{R}^3$. Predict user j 's rating movie i with
 $(\Theta^{(j)})^T x^{(i)}$ stars

$$v(i, j) = 1 \text{ if user } j \text{ has rated movie } i (0 \text{ otherwise})$$

$$y^{(i,j)} = \text{rating by user } j \text{ on movie } i \text{ (if defined)}$$

$\Theta^{(i)}$ = parameter vector for user j

$x^{(i)}$ = feature vector for movie i

For user j , movie i , predicted rating: $(\Theta^{(j)})^T(x^{(i)})$

$m^{(j)}$ = no. of movies rated by user j

To learn $\Theta^{(j)}$:

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

(Optimized)

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

To learn $\Theta^{(1)}, \dots, \Theta^{(n)}$

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

Gradient descent update:

$$\Theta_K^{(j)} := \Theta_K^{(j)} - \alpha \sum_{i=r^{(j)}+1}^n ((\Theta^{(j)})^T x^{(i)} - y^{(i,j)}) x_K^{(i)} \quad (f, k=0)$$

$$\Theta_K^{(j)} := \Theta_K^{(j)} - \alpha \underbrace{\left(\sum_{i=r^{(j)}+1}^n (\Theta^{(j)})^T(x^{(i)}) - y^{(i,j)})^2 + \frac{\lambda}{2} \Theta_K^{(j)} \right)}_{\frac{\partial}{\partial \Theta_K^{(j)}} J(\Theta^{(1)}, \dots, \Theta^{(n)})} \quad (f, K \neq 0)$$

Collaborative Filtering (Basic)

Don't know any features

Given $\Theta^{(1)}, \dots, \Theta^{(n)}$, to learn $x^{(i)}$

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

Given $\theta^{(1)}, \dots, \theta^{(n)}$, to learn $x^{(1)}, \dots, x^{(n)}$.

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

Can guess $\theta \rightarrow x \xrightarrow{\text{new}} \theta \xrightarrow{\text{new}} x \dots$ Converge reasonably

Collaborative Filtering Algorithm

Solve simultaneously

$$J(x^{(1)}, \dots, x^{(n)}, \theta^{(1)}, \dots, \theta^{(n)}) = \frac{1}{2} \sum_{(i,j) \in \mathcal{D}} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{i=1}^m \sum_{k=1}^n (x_k^{(i)})^2 + \frac{\lambda}{2} \sum_{j=1}^n \sum_{k=1}^n (\theta_k^{(j)})^2$$

$$\min_{\substack{x^{(1)}, \dots, x^{(n)} \\ \theta^{(1)}, \dots, \theta^{(n)}}} J(x^{(1)}, \dots, x^{(n)}, \theta^{(1)}, \dots, \theta^{(n)}) \leftarrow \text{minimize simultaneously}$$

$$x \in \mathbb{R}^n$$

$$\theta \in \mathbb{R}^n$$

1. Initialize $x^{(1)}, \dots, x^{(n)}, \theta^{(1)}, \dots, \theta^{(n)}$ to small random values
2. Minimize $J(x^{(1)}, \dots, x^{(n)}, \theta^{(1)}, \dots, \theta^{(n)})$ using gradient descent (or an advanced optimization)
3. For user parameters θ and a movie with $(n \times n)$ features x , predict a star rating $x^{(1)}$

Vectorization: Low Rank Matrix Factorization

Matrix \mathbf{Y} of full ratings

Predicted ratings:

$$[(\theta^{(1)})^T x^{(1)}, \dots, (\theta^{(n)})^T x^{(n)}]$$

$$\begin{bmatrix} (\theta^{(1)})^T x^{(1)} & \dots & (\theta^{(m)})^T x^{(m)} \end{bmatrix}$$

$$x = \begin{bmatrix} -(x^{(1)})^T - \\ \vdots \\ -(x^{(m)})^T - \end{bmatrix} \quad \Theta = \begin{bmatrix} -(\theta^{(1)})^T - \\ \vdots \\ -(\theta^{(m)})^T - \end{bmatrix} \quad X \Theta^T \quad \underset{(i,j) \nearrow}{\begin{pmatrix} (\theta^{(j)})^T x^{(i)} \end{pmatrix}}$$

Finding related movies

For each product i , we learn a feature vector $x^{(i)} \in \mathbb{R}^n$

$x_1 = \text{romance}, x_2 = \text{action}, x_3 = \text{comedy}, x_4 = \dots$

How to find movies j related to movie i ?

$\|x^{(i)} - x^{(j)}\| \leftarrow \text{small} \rightarrow \text{movies } j \text{ and } i \text{ are "similar"}$



Implementation detail: Mean Normalization

User has not rated any movies

$$(\theta^{(i)})^T x^{(i)} = 0 \text{ for all movies}$$

Mean Normalization

$$\text{Matrix } Y \quad \mu = \text{average of rows } Y \rightarrow \overset{\text{normalize}}{Y} = Y - \mu \quad \text{learn } \theta^{(i)}, x^{(i)}$$

For user j , on movie i predict:

$$(\theta^{(j)})^T (\overset{\text{normalize}}{x}^{(i)}) + \mu_i$$