

Advanced Machine Learning

Cristian Perez Jensen

October 27, 2024

Note that these are not the official lecture notes of the course, but only notes written by a student of the course. As such, there might be mistakes. The source code can be found at github.com/cristianpjensen/eth-cs-notes. If you find a mistake, please create an issue or open a pull request.

Contents

<i>1</i>	<i>Paradigms of data science</i>	<i>1</i>
<i>2</i>	<i>Anomaly detection</i>	<i>2</i>

List of symbols

\doteq	Equality by definition
\approx	Approximate equality
\propto	Proportional to
\mathbb{N}	Set of natural numbers
\mathbb{R}	Set of real numbers
$i : j$	Set of natural numbers between i and j . <i>I.e.</i> , $\{i, i+1, \dots, j\}$
$f : A \rightarrow B$	Function f that maps elements of set A to elements of set B
$\mathbb{1}\{\text{predicate}\}$	Indicator function (1 if predicate is true, otherwise 0)
$\boldsymbol{v} \in \mathbb{R}^n$	n -dimensional vector
$\boldsymbol{M} \in \mathbb{R}^{m \times n}$	$m \times n$ matrix
\boldsymbol{M}^\top	Transpose of matrix \boldsymbol{M}
\boldsymbol{M}^{-1}	Inverse of matrix \boldsymbol{M}
$\det(\boldsymbol{M})$	Determinant of \boldsymbol{M}
$\frac{\mathrm{d}}{\mathrm{d}x}f(x)$	Ordinary derivative of $f(x)$ w.r.t. x at point $x \in \mathbb{R}$
$\frac{\partial}{\partial \boldsymbol{x}}f(\boldsymbol{x})$	Partial derivative of $f(\boldsymbol{x})$ w.r.t. \boldsymbol{x} at point $\boldsymbol{x} \in \mathbb{R}^n$
$\nabla_{\boldsymbol{x}}f(\boldsymbol{x}) \in \mathbb{R}^n$	Gradient of $f : \mathbb{R}^n \rightarrow \mathbb{R}$ at point $\boldsymbol{x} \in \mathbb{R}^n$
$\nabla_{\boldsymbol{x}}^2f(\boldsymbol{x}) \in \mathbb{R}^{n \times n}$	Hessian of $f : \mathbb{R}^n \rightarrow \mathbb{R}$ at point $\boldsymbol{x} \in \mathbb{R}^n$

1 Paradigms of data science

Let $\{x_1, \dots, x_n\}$ be i.i.d. samples, generated by an unknown distribution P . Assume that this distribution is in a distribution family,

$$\mathcal{H} = \{p(\cdot \mid \theta) \mid \theta \in \Theta\}.$$

The goal is to learn the parameters θ that fit the data $\{x_1, \dots, x_n\}$ best.

Frequentism. In frequentism, the maximum likelihood estimator (MLE) parameters maximize the following,

$$\begin{aligned} \theta^* &\in \operatorname{argmax}_{\theta \in \Theta} \log p(\{x_1, \dots, x_n\} \mid \theta) \\ &= \operatorname{argmax}_{\theta \in \Theta} \sum_{i=1}^n \log p(x_i \mid \theta). \end{aligned}$$

Bayesianism. Bayesianism assumes that there is a prior over distributions. The maximum a posteriori (MAP) parameters maximize the following,

$$\begin{aligned} \theta^* &\in \operatorname{argmax}_{\theta \in \Theta} \log p(\theta \mid \mathbf{X}) \\ &= \operatorname{argmax}_{\theta \in \Theta} \log p(\{x_1, \dots, x_n\} \mid \theta) \cdot p(\theta) && \text{Bayes' rule.} \\ &= \operatorname{argmax}_{\theta \in \Theta} \log p(\theta) + \sum_{i=1}^n \log p(x_i \mid \theta). \end{aligned}$$

In practice, the prior acts as a regularization term.

Statistical learning. Now, assume that we have labeled samples $\{(x_1, y_1), \dots, (x_n, y_n)\} \subseteq \mathcal{X} \times \mathcal{Y}$, where y is the target variable. Let $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ be a loss function. For a predictor function $f : \mathcal{X} \rightarrow \mathcal{Y}$, we define its risk as the expected loss,

$$\mathcal{R}(f) \doteq \mathbb{E}_{X,Y}[\ell(y, f(x))].$$

In statistical learning, we want to find a function that minimizes the risk. However, since the distribution over X, Y is unknown, we cannot compute $\mathcal{R}(f)$ directly. Instead, we use the empirical risk as a surrogate,

$$\hat{\mathcal{R}}(f) \doteq \frac{1}{n} \sum_{i=1}^n \ell(y_i, f(x_i)).$$

The goal is to obtain the empirical risk minimizer,

$$f^* \in \operatorname{argmin}_{f \in \mathcal{F}} \hat{\mathcal{R}}(f),$$

where \mathcal{F} is a family of functions that we assume f belongs to.

2 Anomaly detection

In anomaly detection, we are given a sample of objects $\mathcal{X} \subseteq \mathbb{R}^d$ with a normal class $\mathcal{N} \in \mathcal{X}$ —the data points that we wish to keep. We wish to construct a function $\phi : \mathcal{X} \rightarrow \{0, 1\}$, such that

$$\phi(x) = 1 \iff x \notin \mathcal{N}.$$

Formally, an anomaly is an unlikely event. Hence, the strategy is to fit a model of a parametric family of distributions to the data \mathcal{X} ,

$$\mathcal{H} = \{p(\cdot | \theta) | \theta \in \Theta\}.$$

Then, we define the anomaly score of x as a low probability $p(x | \theta^*)$ according to the optimal model in this hypothesis class.

Anomaly detection in a high-dimensional space is hard, because the normal class can be very complex. The idea is to project \mathcal{X} down to a lower dimensionality and perform anomaly detection there—hopefully the projected version of the normal class $\Pi(\mathcal{N})$ is less complex. In order to find the optimal linear projection, we will use principal component analysis (PCA).

Furthermore, it has been observed that linear projections of high-dimensional distributions onto low-dimensional spaces resemble Gaussian distributions. Hence, after performing PCA, we will fit a Gaussian mixture model (GMM) to the projected data.

Principal component analysis. The goal of PCA is to linearly project \mathbb{R}^d to \mathbb{R}^{d^-} such that the maximum amount of variance of the data is preserved.¹ Consider the base case $d^- = 1$. Let $\mathbf{u} \in \mathbb{R}^d$ with $\|\mathbf{u}\| = 1$, we project onto \mathbf{u} by inner product,

$$x \mapsto \mathbf{u}^\top x.$$

The sample mean of the reduced dataset is computed by

$$\frac{1}{n} \sum_{i=1}^n \mathbf{u}^\top x_i = \mathbf{u}^\top \bar{x},$$

where \bar{x} is the sample mean of the original dataset. Further, the sample variance of the reduced dataset is

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n (\mathbf{u}^\top x_i - \mathbf{u}^\top \bar{x})^2 &= \frac{1}{n} \sum_{i=1}^n \mathbf{u}^\top (x_i - \bar{x})(x_i - \bar{x})^\top \mathbf{u} \\ &= \mathbf{u}^\top \left(\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^\top \right) \mathbf{u} \\ &= \mathbf{u}^\top \Sigma \mathbf{u}, \end{aligned}$$

where Σ is the covariance matrix of the dataset. Since we want the projection that preserves the maximum variance, we have the following objective,

$$\mathbf{u}^* \in \operatorname{argmax}_{\|\mathbf{u}\|=1} \mathbf{u}^\top \Sigma \mathbf{u}.$$

¹ Components with larger variance are more informative.

The Lagrangian of this problem is

$$\mathcal{L}(\mathbf{u}; \lambda) = \mathbf{u}^\top \Sigma \mathbf{u} + \lambda(1 - \|\mathbf{u}\|^2)$$

with gradient

$$\frac{\partial \mathcal{L}(\mathbf{u}; \lambda)}{\partial \mathbf{u}} = 2\Sigma \mathbf{u} - 2\lambda \mathbf{u} \stackrel{!}{=} 0.$$

So, \mathbf{u} must satisfy $\Sigma \mathbf{u} = \lambda \mathbf{u}$ — \mathbf{u} is an eigenvector of Σ . It is easy to see that this must be the principal eigenvector by rewriting the objective,

$$\begin{aligned} \mathbf{u}^* &\in \operatorname{argmax}_{\|\mathbf{u}\|=1} \mathbf{u}^\top \Sigma \mathbf{u} \\ &= \operatorname{argmax}_{\substack{\|\mathbf{u}\|=1 \\ (\mathbf{u}, \lambda) \in \operatorname{eig}(\Sigma)}} \lambda \|\mathbf{u}\|^2 \\ &= \operatorname{argmax}_{\substack{\mathbf{u} \in \mathbb{R}^d \\ (\mathbf{u}, \lambda) \in \operatorname{eig}(\Sigma)}} \lambda \\ &= \mathbf{u}_1. \end{aligned}$$

For $d^- > 1$, the remaining principal components can be computed with a similar idea. Iteratively, we factor out the previous principal components and do as above on the transformed dataset. For example, to get the second principal component, we first factor out the first principal component,

$$\mathcal{X}_1 \doteq \{\mathbf{x} - \operatorname{proj}_{\mathbf{u}_1}(\mathbf{x}) \mid \mathbf{x} \in \mathcal{X}\} = \{\mathbf{x} - \mathbf{u}_1^\top \mathbf{x} \cdot \mathbf{u}_1 \mid \mathbf{x} \in \mathcal{X}\}.$$

Then, we do the same as above.

Gaussian mixture model. The probability density function (PDF) of a Gaussian mixture model with k components is formalized as a convex combination of Gaussians,

$$p(\mathbf{x}; \boldsymbol{\theta}) = \sum_{j=1}^k \pi_j \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_j, \Sigma_j).$$

The parameters of this model are

$$\boldsymbol{\theta} = \{\pi_j, \boldsymbol{\mu}_j, \Sigma_j \mid j \in [k]\},$$

where $\sum_{j=1}^k \pi_j = 1$ and $\{\Sigma_j \mid j \in [k]\}$ are positive definite. We fit the parameters of this model by maximizing the log-likelihood,

$$\begin{aligned} \boldsymbol{\theta}^* &\in \operatorname{argmax}_{\boldsymbol{\theta} \in \Theta} \log p(\{\mathbf{x}_1, \dots, \mathbf{x}_n\}; \boldsymbol{\theta}) \\ &= \operatorname{argmax}_{\boldsymbol{\theta} \in \Theta} \sum_{i=1}^n \log p(\mathbf{x}_i; \boldsymbol{\theta}) \\ &= \operatorname{argmax}_{\boldsymbol{\theta} \in \Theta} \sum_{i=1}^n \log \sum_{j=1}^k \pi_j \mathcal{N}(\mathbf{x}_i; \boldsymbol{\mu}_j, \Sigma_j). \end{aligned}$$

```

1: Initialize  $\theta_0$ 
2: for  $t \in [T]$  do
3:    $q^* \in \operatorname{argmin}_q E(q, \theta_{t-1})$ 
4:    $\theta_t \in \operatorname{argmax}_\theta M(q^*, \theta)$ 
5: end for
6: return  $\theta_T$ 

```

Note that the above is intractable, so we would like to decompose it into tractable terms that can be computed. Let's assume that we know from which latent variable each data point was generated, then we can compute the MLE of the extended dataset $\{(x_i, z_i) \mid i \in [n]\}$ as

$$\begin{aligned}
\log p(\mathbf{X}, \mathbf{z}; \theta) &= \sum_{i=1}^n \log p(x_i, z_i) \\
&= \sum_{i=1}^n \log(\pi_{z_i} \mathcal{N}(x_i; \mu_{z_i}, \Sigma_{z_i})) \\
&= \sum_{i=1}^n \log \pi_{z_i} + \sum_{i=1}^n \log \mathcal{N}(x_i; \mu_{z_i}, \Sigma_{z_i}),
\end{aligned}$$

which is tractable to maximize. Let q be a distribution over $[k]$, then we can rewrite the log-likelihood into tractable terms,

$$\begin{aligned}
\log p(\mathbf{X}; \theta) &= \mathbb{E}_{z \sim q} [\log p(\mathbf{X}; \theta)] \\
&= \mathbb{E}_{z \sim q} \left[\log \left(\frac{p(\mathbf{X}, z; \theta)}{p(z \mid \mathbf{X}; \theta)} \right) \right] \\
&= \mathbb{E}_{z \sim q} \left[\log \left(\frac{p(\mathbf{X}, z; \theta)}{p(z \mid \mathbf{X}; \theta)} \frac{q(z)}{q(z)} \right) \right] \\
&= \underbrace{\mathbb{E}_{z \sim q} \left[\log \frac{p(\mathbf{X}, z; \theta)}{q(z)} \right]}_{\doteq M(q, \theta)} + \underbrace{\mathbb{E}_{z \sim q} \left[\log \frac{q(z)}{p(z \mid \mathbf{X}; \theta)} \right]}_{\doteq E(q, \theta)}.
\end{aligned}$$

These terms have the following two properties,

$$\begin{aligned}
\log p(\mathbf{X}; \theta) &\geq M(q, \theta), \quad \forall q, \theta \\
\log p(\mathbf{X}; \theta) &= M(q^*, \theta), \quad q^* = p(\cdot \mid \mathbf{X}; \theta), \quad \forall \theta.
\end{aligned}$$

Hence, we can use $M(q^*, \theta)$ as an approximation of $\log p(\mathbf{X}; \theta)$ around θ .

Theorem 2.1 (EM algorithm convergence). Using the expectation-maximization algorithm, $\{\log p(x; \theta_t)\}_{t=0}^T$ is non-decreasing.

Proof. Given a data point x and current parameters θ , we have the following update,

$$\theta' \in \operatorname{argmax}_{\theta \in \Theta} M(q^*, \theta).$$

Algorithm 1. The expectation-maximization algorithm, where

$$\begin{aligned}
M(q, \theta) &\doteq \mathbb{E}_{z \sim q} \left[\log \frac{p(\mathbf{X}, z; \theta)}{q(z)} \right] \\
E(q, \theta) &\doteq \mathbb{E}_{z \sim q} \left[\log \frac{q(z)}{p(z \mid \mathbf{X}; \theta)} \right].
\end{aligned}$$

Hence, we have

$$\log p(\mathbf{x}) = M(q^*, \boldsymbol{\theta}) \leq M(q^*, \boldsymbol{\theta}') \leq \log p(\mathbf{x}; \boldsymbol{\theta}').$$

Thus, $\{\log p(\mathbf{x}; \boldsymbol{\theta}_t)\}_{t=0}^T$ is non-decreasing. ■

Summary. In conclusion, given a set of data points \mathcal{X} with normal points $\mathcal{N} \subseteq \mathcal{X}$, we train an anomaly detector as follows,

1. Fit a projector $\pi : \mathbb{R}^d \rightarrow \mathbb{R}^{d^-}$ using PCA;
2. Fit a probability density function $p(\cdot \mid \boldsymbol{\theta})$ with k components to $\{\pi(\mathbf{x}) \mid \mathbf{x} \in \mathcal{X}\}$ using the EM algorithm;
3. For a point $\mathbf{x} \in \mathcal{X}$, its “anomaly score” is computed by $-\log p(\pi(\mathbf{x}); \boldsymbol{\theta})$.

References