

lesson 8: support vector machines

goal: maximize margin between decision boundary and classes in data

- hard margins: overfits data
- soft margins: allows misclassifications to generalize better with unseen data

$$\text{hinge loss} = \sum_{i=1}^m \max(0, 1 - y^{(i)} (\omega^\top x^{(i)}))$$

cost

$$J(\omega) = C \sum_{i=1}^m \max(0, 1 - y^{(i)} (\omega^\top x^{(i)})) + \frac{1}{2} \sum_{j=1}^n w_j^2$$

$C = \frac{1}{\lambda}$ s.t. λ is the regularization parameter

hypothesis: $h_\omega(x) = \begin{cases} 1 & \text{if } \omega^\top x \geq 0 \\ -1 & \text{if } \omega^\top x < 0 \end{cases}$

goal:

$$\min_w J(\omega) = \min_w C \sum_{i=1}^m \max(0, 1 - y^{(i)} (\omega^\top x^{(i)})) + \frac{1}{2} \sum_{j=1}^n w_j^2$$

Kernels

- during training, features of x are projected to a higher dimensional space
- for a datapoint $x^{(i)}$, the new features $f^{(i)} = \{f_1^{(i)}, \dots, f_m^{(i)}\}$
 - \downarrow n features
 - each datapoint now has m features instead of n
- given that each datapoint x is also a landmark l :

$$f = k(x, l^{(i)}) \quad \text{if features} \rightarrow k(x, l^{(i)}) = 1$$

\downarrow
kernel function | "similarity"

when $x = l^{(i)}$

- updated $J(w) = C \sum_{i=1}^m \max(0, 1 - y^{(i)}(w^T f^{(i)})) + \frac{1}{2} \sum_{j=1}^m w_j^2$

Polynomial Kernel

- $f = k(x, l^{(i)}) = (x \cdot l^{(i)} + r)^d$
 - d \nearrow degree
 - r \downarrow coefficient of the polynomial

radial basis function (RBF) Kernel (aka gaussian kernel)

- $f = k(x, l^{(i)}) = e^{-\gamma(x - l^{(i)})^2}$
- technically, $f \in \mathbb{R}^\infty$, but the "kernel trick" is to calculate f without actually projecting the data point

— when to use SVM vs logistic regression ? —

1 features >> sample size $n \gg m$

\Rightarrow logistic regression or SVM w/out Kernel (linear)

2 features << sample size $n \ll m$

\Rightarrow SVM w/ Kernel (if difference is moderate but still $n \ll m$, feature engineer to ①)

LESSON 9: anomaly detection

$m \times n \rightarrow m$ samples, n features

- given a dataset $X \in \mathbb{R}^n$

- a model $p(x) = \text{probability of } x \text{ being seen in dataset}$

- $h(x) = y = \begin{cases} 1 & \text{if } p(x) < \varepsilon \\ 0 & \text{if } p(x) \geq \varepsilon \end{cases} \xrightarrow{\text{threshold}} \begin{array}{l} \text{anomaly} \\ \Rightarrow \text{normal} \end{array}$

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

\downarrow
jth feature for a
sample $x \in X$

anomaly detection algorithm

- given a datapoint $x \in X$ s.t. $x \in \mathbb{R}^n$:

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

\downarrow mean of each feature j \downarrow variance of feature j

$$\mu = \{\mu_1, \mu_2, \dots, \mu_n\} \quad \sigma^2 = \{\sigma_1^2, \sigma_2^2, \dots, \sigma_n^2\}$$

- $p(x) = \prod_{j=1}^n p(x_j; \mu_j, \sigma_j^2)$ $\xrightarrow{\text{product of } p(x_j) \text{ of features } j}$
 for a single datapoint $x^{(i)}$

3 classify as anomaly if $p(x) < \epsilon$

computing ϵ

includes both normal & anomalous data



1 using the validation dataset, compute $p(x) \forall x$ in dataset

2 extract anomalies in val. data

3 $\epsilon = \max(p(x)) \forall x$ that are anomalies

multivariate same thing but in vector form lol

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

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

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

LESSON 11: PCA & t-SNE

$$\text{Cramer's V} = \sqrt{\frac{\chi^2}{n(k-1)}}$$

sample size \downarrow # of classes \uparrow

\Rightarrow magnitude of association b/w two categorical features

$$\eta^2 = \frac{SS_{\text{between}}}{SS_{\text{total}}} \quad \Rightarrow \text{magnitude of association b/w a categorical & numerical feature}$$

principal component analysis

- a PC explains a % of variability in data

1 center and standardize data

2 calculate covariance $\Sigma = \frac{1}{m} \sum_{i=1}^m (\mathbf{x}^{(i)}) (\mathbf{x}^{(i)})^T$

3 extract eigenvectors & eigenvalues of Σ

4 use the 2 largest eigenvalues & their corresponding eigenvectors
 \downarrow % of variance captured

\downarrow the 2 PC's to use as new axes

5 project original data to PC's

t-SNE

- preserves local relationships between data
- perplexity:
 - determines balance b/w preserving local & global relationships in the data
 - higher P \Rightarrow more separation in low-dimensional space