

## MATHEMATICAL REPRESENTATION:

→ Principal Component Analysis:

Let's consider dataset  $X = \{x_1, x_2, \dots, x_N\}$ , where  $x_n \in \mathbb{R}^D$  with 0 mean and covariance matrix given as,

$$S = \frac{1}{N} \sum_{n=1}^N (x_n - \bar{x})(x_n - \bar{x})^T$$

where  $\bar{x}$  is the mean vector which is zero, hence:

$$S = \frac{1}{N} \sum_{n=1}^N x_n x_n^T$$

Covariance indicates the direction of linear relationship between variables and the matrix

$S$  is the covariance matrix,  $\lambda$  is the eigen value and  $\vec{v}$  is eigen vector. Therefore, principal components?

$$\begin{aligned} S\vec{v} &= \lambda\vec{v} \\ S\vec{v} - \lambda\vec{v} &= 0 \end{aligned}$$

$$\vec{v}^T (S - \lambda I) = 0$$

$(S - \lambda I)$  has to be non-invertible.

$$\det((S - \lambda I)) = 0$$

The eigen values are sorted in decreasing order to form the projection matrix

$$B = [\vec{v}_1 \dots \vec{v}_m] \in \mathbb{R}^{n \times m}$$

for  $m \leq d$ .

Projection matrix thus formed is

$$P = B B^T$$

New feature Subspace  $Z$ :

$$Z = B B^T X$$

$$Z = P X$$

→ Logistic Regression.

Loss Function:

$L(\hat{y}, y)$  = how much  $\hat{y}$  differs from true value of  $y$ .

Let us consider,

$$\hat{y} = P(y=1/n)$$

$\hat{y}_n$  is the probability that  $y=1$  given

$$1 - \hat{y} = P(y = 0 | x)$$

$$P(y|x) = \hat{y}^y (1 - \hat{y})^{1-y}$$

$$\text{if } y = 1 \Rightarrow P(y|x) = \hat{y}$$

Taking the likelihood,  $\log [ \hat{y}^y (1 - \hat{y})^{(1-y)} ]$ .

$$y \log \hat{y} + (1-y) \log (1-\hat{y})$$

$$\Rightarrow -L(\hat{y}, y)$$

$$\Rightarrow \log P(y|x) = -L(\hat{y}, y)$$

The negative log likelihood is to maximize the probability by minimizing loss function.

MLE :

Since we now are using more than two classes, the log of the maximum likelihood function becomes:

$$L(\beta_i^T x) = \sum_{i=1}^N \log (p_i (n_i | \beta_i^T))$$

$$= \sum_{i=1}^N \log \left( \frac{e^{\beta_i^T x_i}}{1 + e^{\beta_i^T x_i}} \right)$$

## The Gradient:

The derivation of the gradient of the maximum likelihood function below:

$$\begin{aligned}\frac{\partial L(\beta, x_i)}{\partial \beta} &= \sum_{i=1}^N \frac{\partial y_i \beta^T x_i}{\partial \beta} - \frac{\partial \ln(1 + e^{\beta^T x_i})}{\partial \beta} \\ &= \sum_{i=1}^N y_i x_i - \frac{x_i e^{\beta^T x_i}}{1 + e^{\beta^T x_i}} \\ &= \sum_{i=1}^N \pi_i (y_i - p_i)\end{aligned}$$

## → Classification

### \* Linear SVM:

SVM hypothesis:

$$h_{w,b}(x) = g(w^T x + b)$$

Class labels are denoted as from 0 to 6;

functional margin and geometric margin,  
optimization problem that the SVM solves

$$\phi(w) = \frac{1}{2} w^T w - \text{minimize}$$

$$\text{Subject to } d_i (w^T x_i + b) \geq 1 \forall i$$

equation for separating hyperplane,  $w$  is the normal to the hyperplane

$$\pi: w^T x^{(i)} + b = 0$$

Functional margin of a hyperplane w.r.t its training example is defined as:

$$\hat{\gamma}^{(i)} = y^{(i)} (w^T x^{(i)} + b)$$

geometric margin of a hyperplane w.r.t the training example,

$$\gamma^{(i)} = \frac{y^{(i)} (w^T x^{(i)} + b)}{\|w\|}$$

This makes our optimization problem as

$$(w^*, b^*) = \underset{w, b}{\operatorname{argmax}} \frac{2}{\|w\|}$$

Such that,

$$y^{(i)} (w^T x^{(i)} + b) \geq 1$$

→ Multinomial Naïve Bayes

The dependent feature vector is  $\{x_1, x_2, \dots, x_n\}$  and the class is  $C_k$ .

$$P(C_k | x_1, \dots, x_n) = \frac{P(C_k) P(x_1, \dots, x_n | C_k)}{P(x_1, \dots, x_n)}$$

Now, for each given class  $C_k$ , each feature vector  $x_i$  is continuously independent of other features.

$$P(x_i | C_k, x_1, \dots, x_n) = P(x_i | C_k)$$

This can be simplified as,

$$P(C_k | x_1, \dots, x_n) = \frac{P(C_k) \prod_{i=1}^n P(x_i | C_k)}{P(x_1, \dots, x_n)}$$

Since  $P(x_1, \dots, x_n)$  is constant,

$$P(C_k | x_1, \dots, x_n) \propto P(C_k) \prod_{i=1}^n P(x_i | C_k)$$

$$\hat{y} = \operatorname{argmax}_k P(C_k) \prod_{i=1}^n P(x_i | C_k)$$

Log probabilities can be used to avoid underflows,

$$\hat{y} = \operatorname{argmax}_k (\ln P(C_k) + \sum_{i=1}^n \ln P(x_i | C_k))$$

The parameter  $\theta_k$  is estimated by maximum likelihood

$$\hat{\theta}_{ki} = \frac{N_{ki} + \alpha}{N_k + \alpha}$$

$$\therefore \hat{y} = \arg \max_k (\ln P(C_k) + \sum_{i=1}^n \ln \frac{N_{ki+g}}{N_k + n})$$