

# 1 Linear Classifier

Feature vectors  $x$ , labels  $y$

$$x \in \mathbb{R}^d$$
$$y \in \{-1, 1\}$$

Training set

$$S_n = \{(x^{(i)}, y^{(i)}), i = 1, \dots, n\}$$

Classifier

$$h: \mathbb{R}^d \rightarrow \{-1, 1\}$$
$$\chi^+ = \{x \in \mathbb{R}^d : h(x) = 1\}$$
$$\chi^- = \{x \in \mathbb{R}^d : h(x) = -1\}$$

Training error

$$\epsilon_n(h) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}\{h(x^{(i)}) \neq y^{(i)}\}$$

Test error (over disjoint set of examples)

$$\epsilon(h)$$

Set of classifiers

$$h \in H$$

**1.1 Linear classifiers through origin**  
Set of all points that satisfies a line through the origin.

$$\theta = \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix}$$
$$X = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

Decision Boundary

$$\{x: \theta_1 x_1 + \theta_2 x_2 = 0\}$$
$$\{x: \theta \cdot X = 0\}$$

Linear Classifier through origin

$$h(x, \theta) = \text{sign}(\theta \cdot X)$$
$$\Theta \in \mathbb{R}^d$$

Decision Boundary

$$\{x: \theta \cdot X + \theta_0 = 0\}$$

Linear Classifier through origin

$$h(x, \theta, \theta_0) = \text{sign}(\theta \cdot X + \theta_0)$$

Linear Classifier through origin

$$h(x, \theta, \theta_0) = \text{sign}(\theta \cdot X + \theta_0)$$

Linear Classifier through origin

$$h(x, \theta, \theta_0) = \text{sign}(\theta \cdot X + \theta_0)$$

Linear Classifier through origin

$$h(x, \theta, \theta_0) = \text{sign}(\theta \cdot X + \theta_0)$$

Linear Classifier through origin

$$h(x, \theta, \theta_0) = \text{sign}(\theta \cdot X + \theta_0)$$

Linear Classifier through origin

$$h(x, \theta, \theta_0) = \text{sign}(\theta \cdot X + \theta_0)$$

Linear Classifier through origin

$$h(x, \theta, \theta_0) = \text{sign}(\theta \cdot X + \theta_0)$$

Linear Classifier through origin

$$h(x, \theta, \theta_0) = \text{sign}(\theta \cdot X + \theta_0)$$

Linear Classifier through origin

$$h(x, \theta, \theta_0) = \text{sign}(\theta \cdot X + \theta_0)$$

Linear Classifier through origin

$$h(x, \theta, \theta_0) = \text{sign}(\theta \cdot X + \theta_0)$$

Linear Classifier through origin

$$h(x, \theta, \theta_0) = \text{sign}(\theta \cdot X + \theta_0)$$

## 1.2 Perceptron through Origin

**Perceptron** $\left(\{(x^{(i)}, y^{(i)}), i = 1, \dots, n\}, T\right)$ :

initialize  $\theta = 0$  (vector);

for  $t = 1, \dots, T$  do

for  $i = 1, \dots, n$  do

if  $y^{(i)}(\theta \cdot x^{(i)}) \leq 0$  then

update  $\theta = \theta + y^{(i)}x^{(i)}$

## 1.5 Perceptron with Offset

**Perceptron** $\left(\{(x^{(i)}, y^{(i)}), i = 1, \dots, n\}, T\right)$ :

initialize  $\theta = 0$  (vector);  $\theta_0 = 0$  (scalar)

for  $t = 1, \dots, T$  do

for  $i = 1, \dots, n$  do

if  $y^{(i)}(\theta \cdot x^{(i)} + \theta_0) \leq 0$  then

update  $\theta = \theta + y^{(i)}x^{(i)}$

update  $\theta_0 = \theta_0 + y^{(i)}$

## 1.6 Margin Boundary

The Margin Boundary is the set of points  $x$  which satisfy  $\theta \cdot x + \theta_0 = \pm 1$ . So, the signed distance from the decision boundary to the margin boundary is  $\frac{1}{\|\theta\|}$ .

$$\frac{y^{(i)}(\theta \cdot x^{(i)} + \theta_0)}{\|\theta\|} = \frac{1}{\|\theta\|}$$

**Hinge Loss (agreement)**

$$\text{Agreement} = z = y^{(i)}(\theta \cdot x^{(i)} + \theta_0)$$
$$\text{Loss}_h(z) = \begin{cases} 0 & \text{if } z \geq 1 \\ 1 - z & \text{if } z < 1 \end{cases}$$

**Regularization** means pushing out the margin boundaries by adding  $\max(\frac{1}{\|\theta\|})$  or  $\min(\frac{1}{2} \|\theta\|^2)$  to the objective function.

## Objective Function

Objective function = average loss + regularization

Objective function is minimized, learning becomes an optimization problem. Using hinge loss and margin boundaries is called **Support Vector Machine** or **Large margin linear classification**:

$$J(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^n \text{Loss}_h(z) + \frac{\lambda}{2} \|\theta\|^2$$

Where  $\lambda > 0$  is called the regularization parameter that regulates how important the margin boundaries are in comparison to the average hinge loss.

## 1.7 Gradient Descent

Assume  $\theta \in \mathbb{R}$  the goal is to find  $\theta$  that minimizes  $J(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^n \text{Loss}_h(y^{(i)}(\theta \cdot x^{(i)} + \theta_0)) + \frac{\lambda}{2} \|\theta\|^2$  through gradient descent.

In other words, we will

- Start  $\theta$  at an arbitrary location:  $\theta \leftarrow \theta_{\text{start}}$
- Update  $\theta$  repeatedly with  $\theta \leftarrow \theta - \eta \frac{\partial J(\theta, \theta_0)}{\partial \theta}$  until  $\theta$  does not change significantly.

Where  $\eta > 0$  is called the stepsize or **learning parameter**.

## 1.8 Stochastic Gradient Descent

$$J(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^n \text{Loss}_h(z) + \frac{\lambda}{2} \|\theta\|^2$$
$$= \frac{1}{n} \sum_{i=1}^n [\text{Loss}_h(z) + \frac{\lambda}{2} \|\theta\|^2]$$

With stochastic gradient descent, we choose  $i \in \{1, \dots, n\}$  at random and update  $\theta$  such that

$$\theta \leftarrow \theta - \eta \nabla_{\theta} [\text{Loss}_h(z) + \frac{\lambda}{2} \|\theta\|^2]$$

## 2 Linear Algebra

### 2.1 Distance

Consider a line  $L$  in  $\mathbb{R}^2$  given by the equation  $L: \theta \cdot x + \theta_0 = 0$  where  $\theta$  is a vector normal to the line  $L$ . Let the point  $P$  be the endpoint of a vector  $x_0$  (so the coordinates of  $P$  equal the components of  $x_0$ ).

The shortest distance  $d$  between the line  $L$  and the point  $P$  is:

$$d = \frac{|\theta \cdot x_0 + \theta_0|}{\|\theta\|}$$

### 3 Regression and classification

Classification:

$$S_n = \{(x^{(t)}, x^{(t)})\} | t = 1, \dots, T$$
$$x^{(t)} \in \mathbb{R}^d, y^{(t)} \in \{-1, 1\}$$

Regression:

$$y^{(t)} \in \mathbb{R}$$
$$f(x, \theta, \theta_0) = \sum_{i=1}^d (\theta_i x_i + \theta_0) = \theta \cdot x + \theta_0$$

### 3.1 Objective for linear regression

The empirical risk  $R_n$  is defined as

$$R_n(\theta) = \frac{1}{n} \sum_{i=1}^n \text{Loss}(y^{(i)} - \theta \cdot x^{(i)})$$

where  $(x^{(t)}, y^{(t)})$  is the  $t$ th training example (and there are  $n$  in total), and Loss is some loss function, such as hinge loss.

Possible to get **closed form solution** for gradient because function is concave. Only possible if the  $d \times d$  matrix  $A$  is invertible. Computationally expensive if dimensions are very high like in bag of words approach.

$$\nabla R_n(\theta) = A\theta - b (=0)$$
$$= A^{-1}b$$

where

$$A = \frac{1}{n} \sum_{i=1}^n x^{(i)}(x^{(i)})^T$$
$$b = \frac{1}{n} \sum_{i=1}^n y^{(i)}x^{(i)}$$

$b$  is a vector with dimensionality  $d$ .

### 3.2 Gradient based Approach

Nudge gradient in the opposite direction to find (local) minima.

$$\nabla_{\theta}(y^{(t)} - \theta x^{(t)})^2 / 2 = (y^{(t)} - \theta x^{(t)}) \nabla_{\theta}(y^{(t)} - \theta x^{(t)}) = -(y^{(t)} - \theta x^{(t)}) \cdot x^{(t)}$$

- initialize  $\theta = 0$
- randomly pick  $t = \{1, \dots, n\}$
- $\theta = \theta + \eta(y^{(t)} - \theta x^{(t)}) \cdot x^{(t)}$

Where  $\eta$  is the learning rate (steps) and the learning rate gets smaller the closer you get

### 3.3 Ridge Regression

Regularization is trying to push away from perfect fit.

$$J_{n,\lambda}(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^n \frac{(y^{(i)} - \theta x^{(i)} - \theta_0)^2}{2} + \frac{\lambda}{2} \|\theta\|^2$$
$$\nabla_{\theta}(J_{n,\lambda}) = \lambda \theta - (y^{(t)} - \theta x^{(t)})x^{(t)}$$

- initialize  $\theta = 0$
- randomly pick  $t = \{1, \dots, n\}$
- $\theta = \theta + \eta \lambda \theta - (y^{(t)} - \theta x^{(t)})x^{(t)} = (1 - \eta \lambda) \theta + \eta(y^{(t)} - \theta x^{(t)})x^{(t)}$

## 3.4 Kernels

$$\phi(x) = [x_1, x_2, x_1^2, \sqrt{2}x_1x_2, x_2^2]$$
$$\phi(x') = [x'_1, x'_2, x'^2_1, \sqrt{2}x'_1x'_2, x'^2_2]$$
$$\phi(x) \cdot \phi(x') = x_1x'_1 + x_2x'_2 + x_1x'^2_1 + 2x_1x'_1x_2x'_2 + x_2x'^2_2$$
$$= (x_1x'_1 + x_2x'_2) + (x_1x'_1 + x_2x'_2)^2$$
$$= (x_1x'_1 + x_2x'_2) + (x_1x'_1 + x_2x'_2)^2$$

### 3.5 Kernel Perceptron

The parameter vector of a perceptron algorithm can also be written as:

$$\theta = \sum_{j=1}^n \alpha_j y^{(j)} \phi(x^{(j)})$$

Where  $\alpha_j$  represents the number of classification mistakes the perceptron algo made. Every time a misclassification happens the parameter vector is updated with the product of the label and the feature vector  $\theta = \theta + y^{(j)}\phi(x^{(j)})$ . The goal of the Kernel Perceptron algo is to find the vector  $\alpha$  with the counts of the misclassifications.

**Kernel Perceptron** $\left(\{(x^{(i)}, y^{(i)}), i = 1, \dots, n, T\}\right)$

initialize  $\alpha_1, \alpha_2, \dots, \alpha_n$ ; to some values

for  $t = 1, \dots, T$  do

for  $i = 1, \dots, n$  do

if  $y^{(i)} \sum_{j=1}^n \alpha_j y^{(j)} K(x^{(i)}, x^{(j)}) \leq 0$  then

update  $\alpha_i = \alpha_i + 1 y^{(i)}$

### 3.6 Radial basis Kernel

$$K(x, x') = e^{-\frac{1}{2} \|x - x'\|^2}$$

## 4 Recommender Systems

### 4.1 K nearest neighbors

The K-Nearest Neighbor method makes use of ratings by  $K$  other “similar users when predicting  $Y_{ai}$ ”. Let  $KNN(a)$  be the set of  $K$  users “similar to user  $a$ ”, and let  $\text{sim}(a, b)$  be a similarity measure between users  $a$  and  $b \in KNN(a)$ . The K-Nearest Neighbor method predicts a ranking  $Y_{ai}$  to be:

$$\widehat{Y}_{ai} = \frac{\sum_{b \in KNN(a)} \text{sim}(a, b) Y_{bi}}{\sum_{b \in KNN(a)} \text{sim}(a, b)}$$

The similarity measure  $\text{sim}(a, b)$  could be any distance function between the feature vectors  $x_a$  and  $x_b$  of users  $a$  and  $b$ , e.g. the euclidean distance  $\|x_a - x_b\|$  and the cosine similarity  $c \cos \theta = \frac{x_a \cdot x_b}{\|x_a\| \|x_b\|}$ .

### 4.2 Collaborative Filtering

Matrix  $Y$  with  $n$  rows (users) and  $m$  columns (Movies) is sparse (entries missing),  $(a, i)$ th entry  $Y_{ai}$  is the rating by user  $a$  of movie  $i$  if this rating has already been given, and blank if not. Goal is to predict matrix  $X$  with no missing entries.

Let  $D$  be the set of all  $(a, i)$  's for which a user rating  $Y_{ai}$  exists, i.e.  $(a, i) \in D$  if and only if the rating of user  $a$  to movie  $i$  exists.

$$J = \sum_{(a,i) \in D} \frac{(Y_{ai} - \frac{(\sum_{j=1}^m U_{Vj}^T)^2}{2})^2}{2} + \frac{\lambda}{2} (\sum_{a,k} U_{ak}^2 + \sum_{i,k} V_{ik}^2)$$
$$u = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix}; v = \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$$
$$v = \begin{bmatrix} 2 \\ 7 \\ 8 \end{bmatrix}$$
$$uv^T = \begin{bmatrix} 2u_1 + 7u_2 + 8u_3 \\ 2u_2 + 7u_3 + 8u_4 \end{bmatrix}$$

Take derivative of Objective function  $J$  with respect to every user, set it to zero and find respective  $u_i$  value:

$$\frac{d}{du_1} \left( \frac{(7-8u_1)^2}{2} + \frac{\lambda}{2} u_1^2 \right) = 0$$
$$\frac{d}{du_2} (J) = 0$$
$$u_1 = \frac{66}{\lambda+68}; u_2 = \frac{16}{\lambda+53}$$

Use resulting values for  $u$  to compute  $uv^T$  compare resulting matrix  $X$  with matrix  $Y$  and start again. Continue until convergence.

### 5 Clustering

Two Views:  
Clustering input:  $S_n = \{x^{(i)} | i = 1, \dots, n\}$   
Clustering output are indexes for the data that partition the data:  $C_1, \dots, C_k$  where  $C_1 \cup C_2 \cup \dots \cup C_k = \{1, 2, \dots, n\}$  and the union of all  $C_i$  's is the original set and the intersection of any  $C_i$  and  $C_j$  is an empty set.  
Representatives of clusters:  $z^{(1)}, \dots, z^{(k)}$ .

Cost of partitioning is the sum of costs of individual clusters:  $\text{cost}(C_1, \dots, C_k) = \sum_{j=1}^k \text{cost}(C_j)$ .

Cost of cluster is sum of distances from data points to the representative of the cluster:  $\text{Cost}(C, z) = \sum_{i \in C} \text{distance}(x^{(i)}, z)$

Cosine similarity:  $\text{cos}(x^{(i)}, x^{(j)}) = \frac{x^{(i)} \cdot x^{(j)}}{\|x^{(i)}\| \|x^{(j)}\|}$  is not sensitive of magnitude of vector (will not react to length).

Euclidean square distance:  $\text{dist}(x^{(i)}, x^{(j)}) = \|x^{(i)} - x^{(j)}\|^2$ . Will react to length.

$$\text{cost}(C_1, \dots, C_k; Z^{(1)}, \dots, Z^{(1)}) = \sum_{j=1}^k \sum_{i \in C_j} \|x^{(i)} - z^{(j)}\|^2$$

### 5.1 The K-Means Algorithm

Only works with Euclidean square distance. Given a set of feature vectors  $S_n = \{x^{(i)} | i = 1, \dots, n\}$  and the number of clusters  $K$  we can find cluster assignments  $C_1, \dots, C_K$  and the representatives of each of the  $K$  clusters  $z_1, \dots, z_K$ :

- Randomly select  $z_1, \dots, z_K$
- Iterate
  - Given  $z_1, \dots, z_K$ , assign each data point  $x^{(i)}$  to the closest  $z_j$ , so that  $\text{Cost}(z_1, \dots, z_K) = \sum_{i=1}^n \min_{j=1, \dots, K} \|x^{(i)} - z_j\|^2$
  - Given  $C_1, \dots, C_K$  find the best representatives  $z_1, \dots, z_K$ , i.e. find  $z_1, \dots, z_K$  such that  $z_j = \arg \min_z \sum_{i \in C_j} \|x^{(i)} - z\|^2$

The best representative is found by optimization (gradient with respect to  $z^{(j)}$ , setting to zero and solving for  $z^{(j)}$ ).

It is the centroid of the cluster:  $z^{(j)} = \frac{\sum_{i \in C_j} x^{(i)}}{|C_j|}$

The clustering output that the K-Means algorithm converges to depends on the initialization.

### 5.2 K-Medoids Algorithm

Finds the cost-minimizing representatives  $z_1, \dots, z_K$  for any distance measure. Uses real data points for initialization.

- Randomly select  $\{z_1, \dots, z_K\} \subseteq \{x_1, \dots, x_n\}$
- Iterate
  - Given  $z_1, \dots, z_K$ , assign each data point  $x^{(i)}$  to the closest  $z_j$ , so that  $\text{Cost}(z_1, \dots, z_K) = \sum_{i=1}^n \min_{j=1, \dots, K} \|x^{(i)} - z_j\|^2$

(b) Given  $C_j \in \{C_1, \dots, C_K\}$  find the best representative  $z_j \in \{x_1, \dots, x_n\}$  such that  $\sum_{x^{(i)} \in C_j} \text{dist}(x^{(i)}, z_j)$  is minimal

## 6 Generative Models

Understand structure of data probabilistically.

### 6.1 Multinomial Models

Fixed Vocabulary  $W$

Multinomial model  $M$  to generate text in documents.

Document  $D$

Likelihood of generating certain word  $w \in W$ :  $p(w|\theta) = \theta_w$  where  $\theta_w \geq 0$  and  $\sum_{w \in W} \theta_w = 1$ .  
Likelihood function:

$$P(D|\theta) = \prod_{i=1}^n \theta_{w_i}$$
$$= \prod_{w \in W} \theta_w^{\text{count}(w)}$$

Toy Example:

$$\theta_1: \theta_{\text{cat}} = 0.3; \theta_{\text{dog}} = 0.7$$
$$\theta_2: \theta_{\text{cat}} = 0.9; \theta_{\text{dog}} = 0.1$$
$$D = [\text{cat}, \text{cat}, \text{dog}]$$
$$P(D|\theta_1) = 0.3^2 \cdot 0.7 = 0.063$$
$$P(D|\theta_2) = 0.9^2 \cdot 0.1 = 0.081$$

Maximum likelihood:

$$\max_{\theta} P(D|\theta) = \max_{\theta} \prod_{w \in W} \theta_w^{\text{count}(w)}$$
$$\log \prod_{i=1}^n \theta_w^{\text{count}(w)} = \sum_{w \in W} \text{count}(w) \log(\theta_w)$$
$$W = \{0, 1\}; \theta_0 = 0; \theta_1 = (1 - \theta)$$
$$\frac{d}{d\theta} (\text{count}(0) \log(\theta) + \text{count}(1) \log(1 - \theta)) = \frac{\text{count}(0)}{\theta} - \frac{\text{count}(1)}{1 - \theta} = 0$$
$$\hat{\theta} = \frac{\text{count}(0)}{\text{count}(0) + \text{count}(1)}$$

For any length of  $W$ :

$$\hat{\theta} = \frac{\text{count}(w)}{\sum_{w' \in W} \text{count}(w')}$$

### 6.2 Prediction

Goal: categorize between minus and plus class. Both classes have a associated parameter  $\theta^+$  and  $\theta^-$   
Class conditional distribution:

$$\log \left( \frac{P(D|\theta^+)}{P(D|\theta^-)} \right) = \begin{cases} \geq 0, & + \\ < 0, & - \end{cases}$$

Model is the same as a linear classifier through origin:

$$\log(P(D|\theta^+)) - \log(P(D|\theta^-)) = \log \prod_{w \in W} \theta_w^{+\text{count}(w)} - \log \prod_{w \in W} \theta_w^{-\text{count}(w)} = \sum_{w \in W} \text{count}(w) \log(\theta_w^{+\text{count}(w)}) - \sum_{w \in W} \text{count}(w) \log(\theta_w^{-\text{count}(w)}) = \sum_{w \in W} \text{count}(w) \log \frac{\theta_w^{+\text{count}(w)}}{\theta_w^{-\text{count}(w)}} = \sum_{w \in W} \text{count}(w) \bar{\theta}_w$$

