


## Linear Predictive Model

$$\hat{y}_i = b^T X_i + b$$

$$P(y_i=1|x_i) = \frac{\exp(\hat{y}_i)}{1 + \exp(\hat{y}_i)}$$

Sigmoid function 

## Softmax Function

$$\text{softmax}(\hat{y}_i)_j = \frac{\exp(\hat{y}_{ij})}{\sum_{l=1}^J \exp(\hat{y}_{il})}$$

$$P(y_i=j|x_i) = \text{softmax}(\hat{y}_i)_j$$

## Empirical Risk Minimization

$$b^* = \underset{b}{\text{argmin}} \frac{1}{N} \sum_{i=1}^N \ell(y_i, f(x_i; b))$$

## Loss Function

$$\ell(y, \hat{y}) = -y \log \hat{y} - (1-y) \log(1-\hat{y}) \quad [0, 1]$$

$$\ell(y, \hat{y}) = -\sum_i y_i \log(\hat{y}_i) \quad [0, 1, 2, \dots]$$

## K-Nearest Neighbors

- classifier:

1. set  $K=1, 2, \dots$
2. get  $X_i \in X_0, X_1, X_2, \dots$  distance
3. select  $K$  nearest neighbor
4. majority voting

- regressor:

1. set  $K$
2. get  $X_i \in X_0, X_1, X_2, \dots$  distance
3.  $\hat{y} = \frac{1}{K} \sum_{i \in K_{\text{nn}}} y_i$

or

$$\hat{y} = \frac{\sum_i \frac{1}{d(x, x_i)} y_i}{\sum_i \frac{1}{d(x, x_i)}}$$

Overfitting  $\leftrightarrow$  underfitting

Train  $\rightarrow$  Valid  $\rightarrow$  Test

## Cross Validation

- reduce variance of performance estimates
- $N$ -fold cross validations

## Metric

$$\text{MSE} = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

Accuracy:

$$\frac{\# \text{ correct guesses}}{\# \text{ of data points}}$$

99% sensitive

$\hookrightarrow$  99% of (+) cases  $\rightarrow$  predict (+)

99% specific

$\hookrightarrow$  99% of (-) cases  $\rightarrow$  predict (-)

(-) predict (+) predict

(+) true FN TP  
(-) true TN FP

$$\text{accuracy} = \frac{TN + TP}{\text{All}}$$

## K-fold cross validation

San fold in folds:

valid  $\leftarrow$  fold  
train  $\leftarrow$  folds \setminus fold

San  $k$  in  $k$ -list

train on train  
valid on valid  
get - performance  
store performance

evaluate score on each  $k$

$$k^* \leftarrow \underset{k}{\text{argmax}} \text{Score}_{\text{avg}}(k)$$

Feature selection ANOVA analysis of var  
univariate explained var =  $\frac{N_0(\bar{X}_0 - \bar{X})^2 + N_1(\bar{X}_1 - \bar{X})^2}{(N_0-1)\sigma_0^2 + (N_1-1)\sigma_1^2}$   
F-score =  $\frac{\text{explained var}}{\text{unexplained var}}$

How different the group means are  
 $\rightarrow$  random noise

high: group means different; feature rich  
low: " close; feature poor

Sequential feature selection (forward)

1. no feature
2. while model improve
  - add feature  $f_i$  in  $f$ -list
  - evaluate
  - choose best feature  $f_i$
  - add  $f_i^*$

univariate  $v$

PRO fast  
easy  
interpretable

sequential  
capture interaction  
accuracy  $\uparrow$

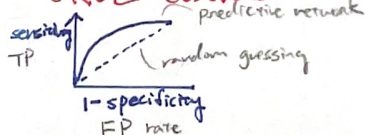
CON ignore correlation

greedy  
slow  
oversit

Sequential feature selection (backward)

1. all features
2. while model improve
  - remove feature  $f_i$  in  $f$ -list
  - remove  $f_i$
  - $\rightarrow$  improves model

## Receiving Operating Characteristic (ROC curve)



AUC: integral of the plot under the curve

1. classifier output  $p \in [0, 1]$
2. vary threshold  $0 \rightarrow 1$   
 $\rightarrow$  let  $p > \tau = (+)$
3. for each  $\tau$ ,  
get TPR, FPR
4. plot TPR  $v$  FPR

## Confusion Matrix

Predict	0	1	...	N
Actual	0	1	...	N
	TP	FN		
	FP	TN		

$$\text{TPR} = \frac{TP}{TP + FN}$$

$$\text{FPR} = \frac{FP}{FP + TN}$$

## SVM



$$1. w^T X + b \geq m \quad (o)$$

$$w^T X + b \leq -m \quad (x)$$

$$w^*, b^* = \underset{w, b}{\text{argmax}} m \quad \text{s.t. } \|w\|_2 = 1$$

$$2. w^T x_i + b \geq m$$

$$\rightarrow \frac{w^T}{m} x_i + \frac{b}{m} \geq 1$$

$$3. w^* = \underset{w}{\text{argmin}} \|w\|_2^2$$

$$w^T x + b > 1 \quad \forall (o)$$

$$w^T x + b < -1 \quad \forall (x)$$

slack variables

$$w^* = \underset{w}{\text{argmin}} \|w\|_2^2 + \lambda \sum_{i=1}^n \psi_i$$



## Kernel

$$x \rightarrow \phi(x)$$

linear  
polynomial  
exponential  
radial basis

## decision tree

cross-entropy loss

$$\text{Lm} = - \sum_{k=1}^K p_{mk} \log p_{mk}$$

$$p_{mk} = \frac{\# \text{ samples w/ outcome } k \text{ in node } m}{\# \text{ of samples in node } m}$$

random forest aiding decision tree

- smoother boundaries
- ensemble
- reduce variance
- increase bias

## K-means

$$0. \text{SSE} = \sum_{i=1}^n \|X_i - C_g\|_2^2$$

1. given  $x_1, x_2, \dots, x_n \in \mathbb{R}^d$   
cluster no.  $k$

2. initialized  $C_0, C_1, C_2, \dots, C_k$

3. repeat until convergence

• assignment:

for each data  $i$ , do:

$$g_i \leftarrow \underset{j \in \{1, \dots, k\}}{\text{argmin}} \|X_i - C_j\|_2^2$$

$\hookrightarrow$  index of nearest centroid

• centroid update

$$C_j = \frac{1}{n_j} \sum_{i \in S_j} X_i$$

$p$  is the  $p$ th feature

Standard Scaling (Z-scoring)

$$\tilde{x}_{ip} = \frac{x_{ip} - \bar{x}_p}{\hat{\sigma}_p}$$

$$\tilde{x}_p = 0$$

$$\tilde{\sigma}_p = 1$$

Min-Max Rescaling

$$\tilde{x}_{ip} = \frac{x_{ip} - \bar{x}_p^{\min}}{\bar{x}_p^{\max} - \bar{x}_p^{\min}}$$

Max-abs Rescaling

$$\tilde{x}_{ip} = \frac{x_{ip}}{\bar{x}_p^{\max}}$$

LASSO  $v$  RIDGE

- many  $w_i$  - smoother  
 $\rightarrow 0$

- L1-norm - L2-norm

- some  $f_i$  - most  $f_i$   
irrelevant useful  
use this

$\checkmark$   
combine  
 $\parallel$

Elastic Net



## GMM

$$P(x) = \sum_{k=1}^K \pi_k \mathcal{N}(x | \mu_k, \Sigma_k)$$

$$\pi_k \Rightarrow \sum_{k=1}^K \pi_k = 1$$

## Expectation - Maximization Algo

$$X = \{x_1, x_2, \dots, x_N\}$$

→ assume this is generated w/ a Gaussian Mixture Model

### 1. hidden variables $r_{ik}$

→ how much does  $k^{\text{th}}$  distribution "own" datapoint  $x_i$ ?

$$r_{ik} = P(z_i = k | x_i)$$

### 2. initial guess

$$\mu_k, \Sigma_k, \pi_k$$

### 3. Expectation step

$$r_{ik} = \frac{\pi_k \mathcal{N}(x_i | \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(x_i | \mu_j, \Sigma_j)}$$

→ compute the relative prob. that  $x_i$  belongs to  $k^{\text{th}}$  distribution

### 4. Maximization step

$$N_k = \sum_{i=1}^N r_{ik}$$

$$\pi_k = N_k / N$$

$$\mu_k = \frac{1}{N_k} \sum_{i=1}^N r_{ik} x_i$$

$$\Sigma_k = \frac{1}{N_k} \sum_{i=1}^N r_{ik} (x_i - \mu_k)(x_i - \mu_k)^T$$

$$5. \mathcal{L} = \sum_{i=1}^N \log \left( \sum_{k=1}^K \pi_k \mathcal{N}(x_i | \mu_k, \Sigma_k) \right)$$

6. loop → until  $\mathcal{L}$  not improving

## PCA

w/  $n$  features  
in data pts

$$0. X = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1n} \\ x_{21} & x_{22} & \dots & x_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{nn} \end{bmatrix} \in \mathbb{R}^{n \times n}$$

$k = \begin{bmatrix} n \\ \vdots \\ 1 \end{bmatrix}$

$$1. Z = \frac{1}{n-1} X^T X$$

$$Z \in \mathbb{R}^{n \times n}$$

$$2. \sum \lambda_i = \text{tr}(Z) \in \mathbb{R}^n$$

cov. matrix  
vec of  
"principal components"

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$$

$$3. W = \begin{bmatrix} | & | & \dots & | \\ v_1 & v_2 & \dots & v_n \\ | & | & \dots & | \end{bmatrix} \in \mathbb{R}^{n \times l}$$

$$l \leq n$$

$$4. Z = X \cdot W \in \mathbb{R}^{n \times l}$$

$\mathbb{R}^{n \times n} \mathbb{R}^{n \times l}$

$$5. \hat{X} = Z W^T$$

## t-SNE

+ distributed stochastic neighbor embedding

$$X = \{x_1, x_2, \dots, x_n\}$$

perplexity

local size

K

target dim

lr, max-iter

1. for  $x_i$  in  $X$ :

find  $\pi_i$  s.t. entropy of  $\pi_i$  = perplexity

for  $x_j \neq x_i$ :

$$P_{ij} = \exp(-\|x_i - x_j\|^2 / (2\sigma_i^2))$$

$$P_{ij} = P_{ij} / \sum_j P_{ij}$$