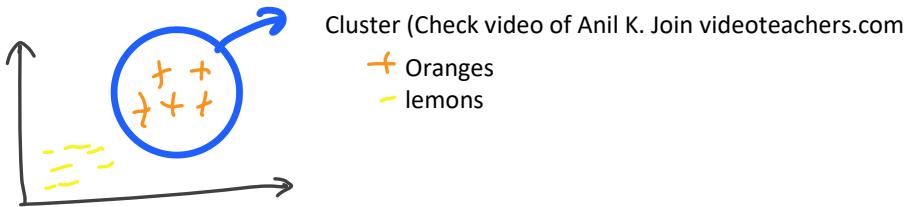


Intro

8 Şubat 2017 Çarşamba 13:45

Online Courses : Yaser Abou Moustafa - Caltech
Nando Freitas - Oxford

Book : Ffont.Intro to Machine Learning 3rd -Ethem Alpaydin
Machine Learning & Pattern Recognition- Christopher Bishop

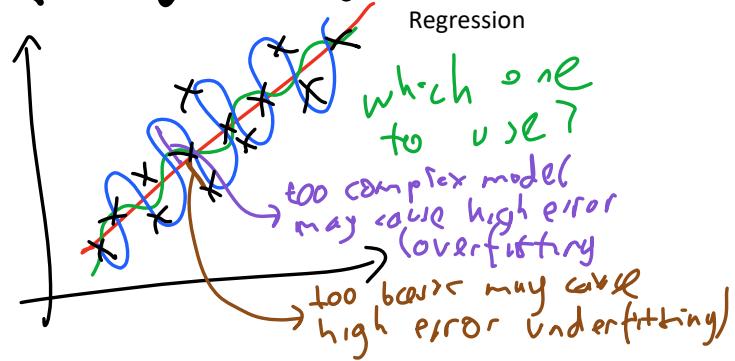
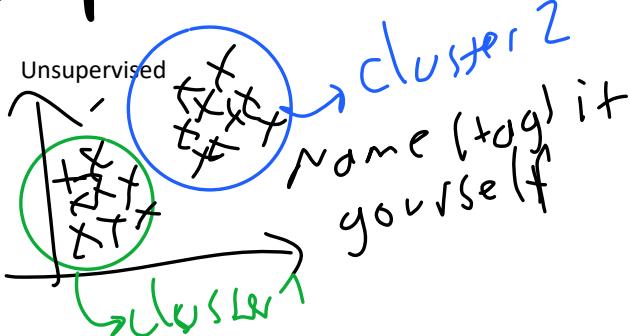


Data is given : Supervised Learning -find regulars of data points

This graph two dimensional but can be increased Like color and sugar Level

We need to turn color into math function so we test our features with training dataset. And we have to do it more than once to increase statistical accuracy, least classification error and highest accuracy with biggest data set is the winner. We will use this dividing rule in the application

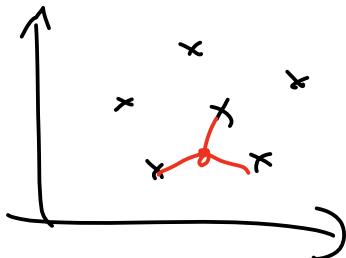
But also which features to use is important too. (Feature engineering)



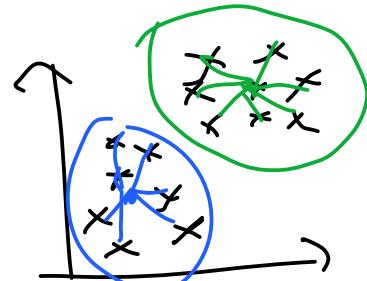
Classifiers

8 Şubat 2017 Çarşamba 14:54

First clustering K-NN



Find three closest points and assume new point is their type
Drawbacks:
KEEPING ALL TRAINING DATA POINTS
COMPUTING ALL THE DISTANCES



K-means clustering

K=2

Iterate to cluster centers

You may come up with Totally different clusters

Need to know numbers Of clusters

Group and calculate Mean until distances does not change very much

Other statistics like Hubert's (not covered in class)

You can also calculate distance between cluster means and merge if its lower than threshold

- ① Supervised Learning → *Decision boundary*
-
- $x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$
- +,- labels KNN algoritm
- ② Unsupervised Learning (Clustering)
-
- K-medians
- ③ Regression
-
- Keep the model simple as possible

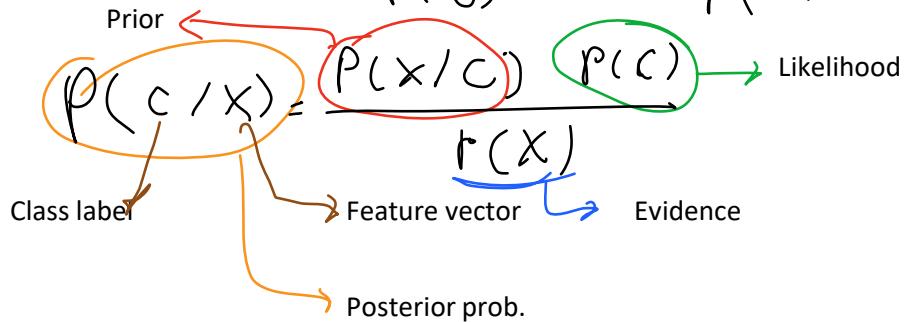
You are given the data. You will select data randomly or Smartly (later in course) and create a model
If model has low error for both training and real data model is good

But it's impossible to get right data with random choosed training data (possibility might be used to select right data- check statistics notes)

Statistics

15 Şubat 2017 Çarşamba 14:07

$$P(A|B) = \frac{P(A|B)}{P(B)} = \frac{P(B|A)P(A)}{P(B)}$$



Chalkboard content:

$$P(C=1|M_1=0, M_2=1) = \frac{P(M_1=0, M_2=1|C=1)P(C=1)}{P(M_1=0, M_2=1)}$$

$$= \frac{4}{13} = \frac{4}{20}$$

$$P(C=0|M_1=0, M_2=1) = \frac{P(M_1=0, M_2=1|C=0)P(C=0)}{P(M_1=0, M_2=1)}$$

$$= \frac{13}{20}$$

Distribution & Expectations-Barneuli Random

A binary random variable X is a mapping from the sample space Ω to a discrete space $E = \{0, 1\}$

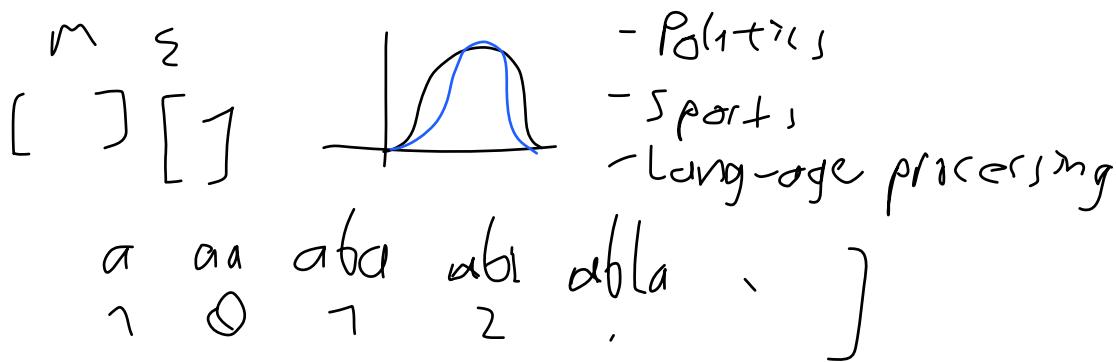
$$X(\omega) \sim \rightarrow E$$

↳ set of ω

Probability Distribution (with random variable - check statistic notes)

$$P(X=\omega) = \underbrace{\varphi^{\omega}}_{\text{Parameter}} \otimes (1-\varphi)^{1-\omega}$$





Naïve Bayes Classifiers : Features Conditionally Independent

$$P(x_1, x_2, \dots, x_n | C) \xrightarrow{\text{Assume independent}}$$

$$P(x_j | C) P(x_i | C)$$

$$P(C | x) = \frac{P(x | C) P(C)}{P(x)}$$

$$P(C | x_1, x_2, \dots, x_n) = \frac{P(x_1, x_2, \dots, x_n | C) P(C)}{P(x)} \Rightarrow \text{Naïve Bayes}$$

$$= \frac{\prod_{i=1}^n P(x_i | C) P(C)}{P(x)} \quad \text{Select } C \text{ that maximizes } P(C | x)$$

Naïve Bayes works as

k classes d dimensional features

$$P(C_i) = \frac{\# \text{ that belongs to } C_i}{\# T} \xrightarrow{\text{Size of training set}}$$

x are contours

$$P(x_i | C_i) \approx \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(\frac{(x - \mu)^2}{2\sigma^2}\right)$$

Let features be X_i

If the features are categorized and if the feature X_i has category j (for binary features; 0, 1)

$$P(X_i=j | C_i) = \frac{\# j}{\text{Total number of examples in class } C_i}$$

Method is same with bayes but assuming possibilities are independent, calculating mean and covariance matrix becomes unnecessary therefore, number of calculations will be reduced.

Age	25	30	35	40	45	50	55	60	65	70	75	80
Income	low	low	low	low	low	high						
Student	yes	yes	yes	yes	yes	yes	yes	yes	yes	yes	yes	yes
Credit	bad	bad	bad	bad	bad	bad	bad	bad	bad	good	good	good
Buy	no	no	no	no	no	no	no	no	no	no	yes	yes

Classify $X = \{ \text{age} = \text{youth}, \text{income} = \text{medium}, \text{student} = \text{yes}, \text{credit} = \text{bad} \}$

$$P(C_2) = 5/14 \quad P(C_1) = 9/14$$

$$P(\text{age} = \text{youth} | C_2) = 2/9$$

$$P(\text{age} = \text{youth} | C_1) = 3/5$$

$$P(\text{student} = \text{yes} | C_1) = 6/9$$

$$P(\text{student} = \text{yes} | C_2) = 2/5$$

$$P(\text{income} = \text{m} | C_1) = 5/9$$

$$P(\text{income} = \text{m} | C_2) = 2/5$$

$$P(\text{credit} = \text{f} | C_1) = 0$$

$$P(\text{credit} = \text{f} | C_2) = 2/5$$

$$P(X | C_1) = \frac{2}{7} \cdot \frac{6}{9} \cdot \frac{5}{9} = 0.049$$

What if special set does not exist $P(X)$ will be zero. This set may happen in real life, but it might included in training set, it may be caused by little training data. How can we solve this problem without expanding dataset available.

$$P(X) = 0 \Rightarrow \frac{0+1}{9+3} \text{ number of different features}$$

0 feature { If there is a lot of
kafunki } features the will
deger datibildi { be negligible

Linear Prediction

8 Mart 2017 Çarşamba 15:03

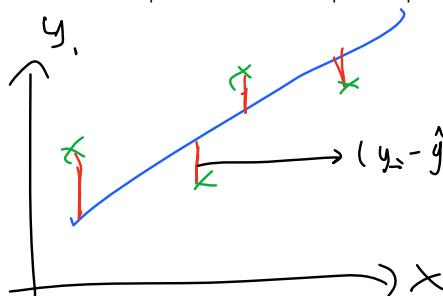
Linear Regression (Least Squares Model)
Given a training set n instances $\{X_{1in}, Y_{1in}\}$

Each input $X_i \in \mathbb{R}^{1 \times d}$ a vector with d attributes (features)

Output y_i referred as target $y_i \in \mathbb{R}$ $n=4$

Wind Speed	#people	Energy Req.
100	2	5
50	20	18
45	44	21
60	35	40

Learn a model that represents the effect of inputs to outputs



$$f(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 \rightarrow \text{local}$$

$$\sum_{i=1}^n (y_i - \hat{y}_i)^2 = J(\theta)$$

$$\sum_{i=1}^n (y_i - (\theta_0 + \theta_1 x_i))^2$$

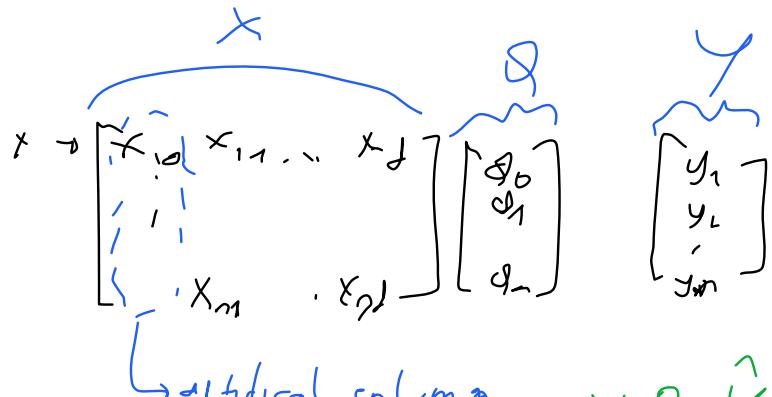
$$\underset{\theta}{\operatorname{arg\,min}} \sum_{i=1}^n (y_i - \theta_0 - \theta_1 x_i)^2$$

$$\frac{\partial J(\theta)}{\partial \theta_0} = 0$$

$$\frac{\partial J(\theta)}{\partial \theta_1} = 0$$

$$\frac{\partial J(\theta)}{\partial \theta_d} = 0$$

$$y_i = \theta_0 + \theta_1 x_1 + \theta_d x_d \Rightarrow \hat{y} = \sum_{j=1}^d \theta_j x_j + \theta_0 = \sum_{j=1}^d x_j \theta_j$$



$$(Y - \hat{Y})^\top (Y - \hat{Y}) = J(\theta) \Rightarrow \text{make square differences in matrix to make working with higher dimensions easier}$$

$$[y_1 - \hat{y}_1, y_2 - \hat{y}_2] \begin{bmatrix} y_1 - \hat{y}_1 \\ y_2 - \hat{y}_2 \end{bmatrix} = (y_1 - \hat{y}_1)^2 + (y_2 - \hat{y}_2)^2$$

$$L' \rightarrow Y_1 - \hat{Y}_1 + Y_2 - \hat{Y}_2 \rightarrow \begin{pmatrix} -1 & 1 \\ 1 & -1 \end{pmatrix} = (Y_1 - \hat{Y}_1) + (Y_2 - \hat{Y}_2)$$

$$(Y - X\theta)^T (Y - X\theta) = J(\theta)$$

$$= Y^T Y - Y^T X \theta - (X\theta)^T + (X\theta)^T (X\theta)$$

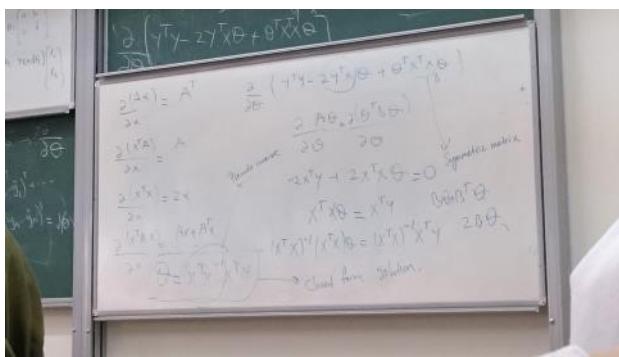
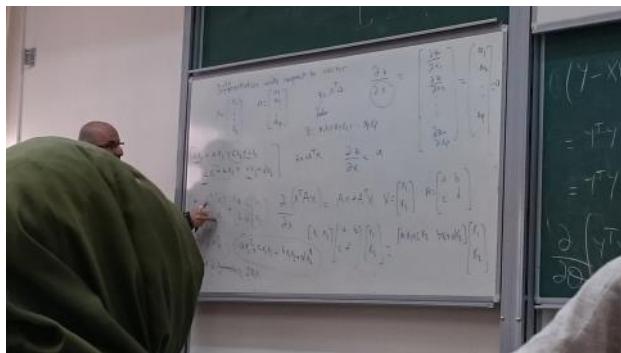
$$= Y^T Y - Y^T X \theta - \theta^T X^T Y + \theta^T X^T X \theta$$

$$\frac{\partial}{\partial \theta} (Y^T Y - 2Y^T X \theta + \theta^T X^T X \theta) \quad d = d^T$$

Differentiation w.r.t. θ gives the gradient

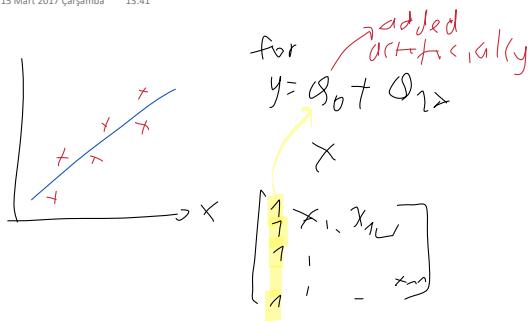
$$x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \quad d = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix} \quad z = x^T d$$

$$z = a_1 x_1 + a_2 x_2 + \dots + a_n x_n \quad t = a^T x \quad \frac{\partial z}{\partial x} = \begin{bmatrix} \frac{\partial z}{\partial x_1} \\ \frac{\partial z}{\partial x_2} \\ \vdots \\ \frac{\partial z}{\partial x_n} \end{bmatrix}$$



Complexity high because of large matrix multiplications and matrix derivatives

θ 's dimension $(d+1) \times 1$ - All θ 's found



Require the inversion of:

$X^T X \rightarrow$ If system eq. are poorly conditioned you may not be able to obtain the inverse

A solution to this problem is to add a small element (ϵ^{12}) to the diagonal:

$$\hat{\theta} = (X^T X + \epsilon^2 I_{(d+1)(d+1)})^{-1} X^T Y$$

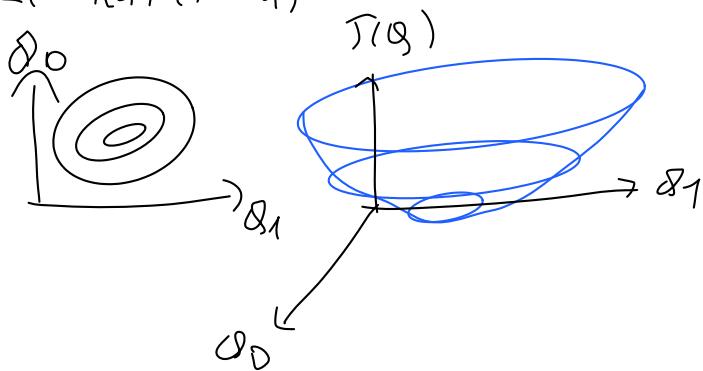
$$X \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_d \end{bmatrix} \begin{bmatrix} \bar{y}_1 \\ y_2 \\ \vdots \\ \bar{y}_n \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

$$J(\theta) = [y_1 - \bar{y}_1 \quad y_2 - \bar{y}_2 \quad \dots \quad y_n - \bar{y}_n]^T \begin{bmatrix} y_1 - \bar{y}_1 \\ y_2 - \bar{y}_2 \\ \vdots \\ y_n - \bar{y}_n \end{bmatrix} = (y - \bar{y})^T (y - \bar{y})$$

error function

$$J(\theta) = (\bar{y} - y)^T (\bar{y} - y) = \sum_{i=1}^n (y_i - \theta_0 - \theta_1 x_1 - \theta_2 x_2 - \dots - \theta_d x_d)^2$$

$$J(\theta) = (y - X\theta)^T (y - X\theta)$$



To optimize (in machine learning and neural networks) we must minimize $J(\theta)$. Now we will learn about iterative methods to minimize $J(\theta)$

Don't let big θ 's to dominate others

$$J(\theta) = (y - X\theta)^T (y - X\theta) + \frac{\lambda}{2} \theta^T \theta$$

↳ Regularized cost function ↳ Regularization ↳ Regularization parameter

Taking derivative of $J(\theta)$

Reminder

$$\frac{\partial^T Y}{\partial X} - \alpha$$

$$\frac{\partial X^T A X}{\partial X} = A^T Y + A X$$

Taking derivative of $J(\theta)$

$$S(\theta) = Y^T Y - Y^T X \theta - \theta^T X^T Y + \theta^T X^T X \theta + \delta^2 \theta^T \theta$$

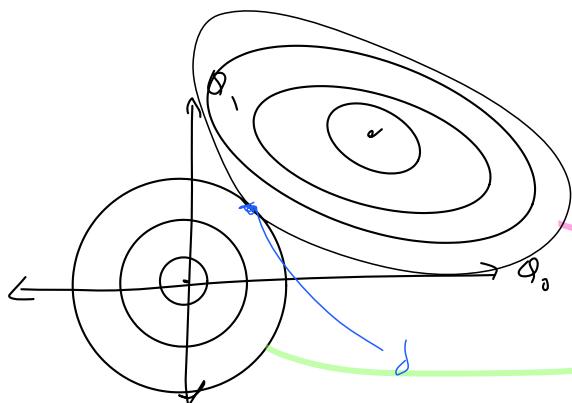
$$\frac{\partial}{\partial \theta} [Y^T Y - 2Y^T X \theta + \theta^T X^T Y + \delta^2 \theta^T \theta] \\ - 2X^T Y + 2X^T X \theta + 2\delta^2 \theta = 0$$

$$X^T X \theta + \delta^2 \theta = X^T Y$$

$$(X^T X + \delta^2 I) \theta = X^T Y$$

$$\theta = (X^T X + \delta^2 I)^{-1} X^T Y$$

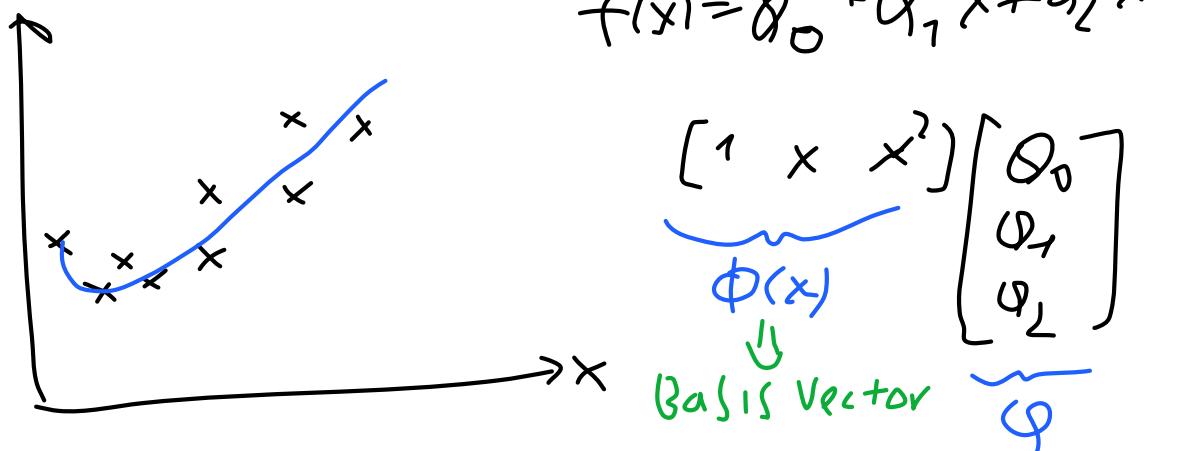
the parameter must be fine tuned from the dataset



$$J(\theta) = (Y - X\theta)^T (Y - X\theta) + \delta^2 \theta^T \theta$$

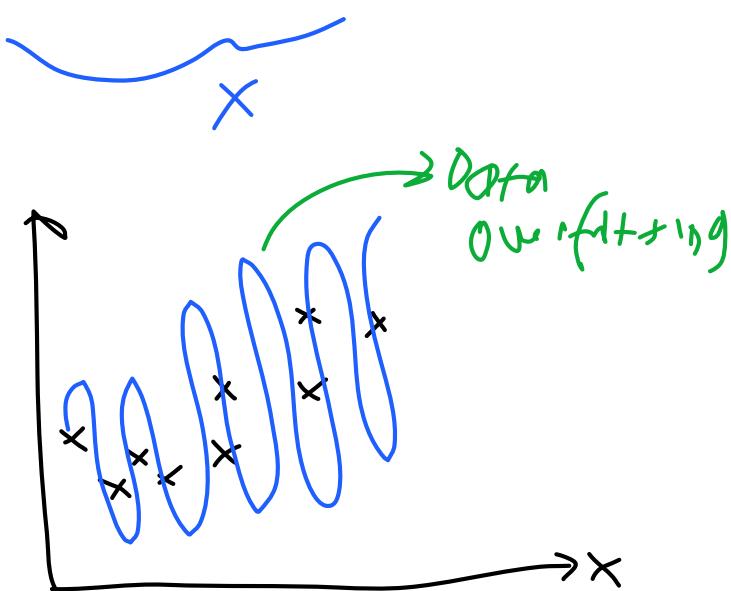
NonLinear Prediction

15 Mart 2017 Çarşamba 14:42

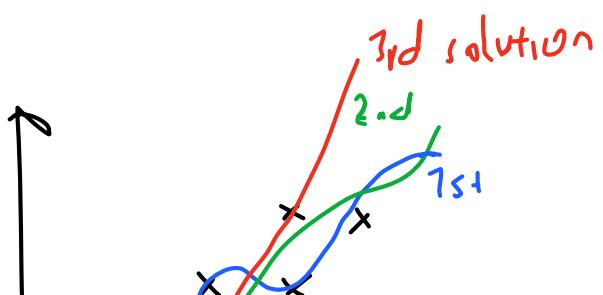


Assume x^2 is another function and create basis and q matrix according to that.

$$\begin{bmatrix} 1 & x_1 & x_1^2 & \dots & x_L & x_L^2 & \dots & x_n & x_n^2 \\ 1 & x_1 & x_1^2 & \dots & x_L & x_L^2 & \dots & x_n & x_n^2 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 1 & x_1 & x_1^2 & \dots & x_L & x_L^2 & \dots & x_n & x_n^2 \end{bmatrix} \boldsymbol{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_{19} \end{bmatrix}$$

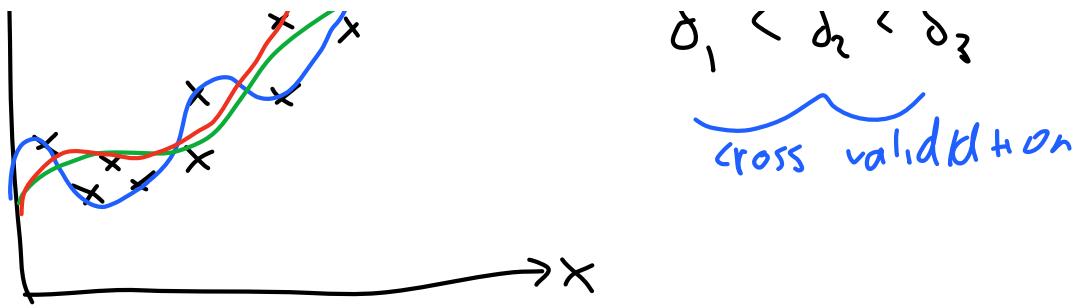


If you don't Use regulator yor. Data will be overfitted



$$f(x) = \theta_0 + \theta_1 x + \theta_4 x^4$$

$$\delta_1^2 < \delta_2^2 < \delta_3^2$$



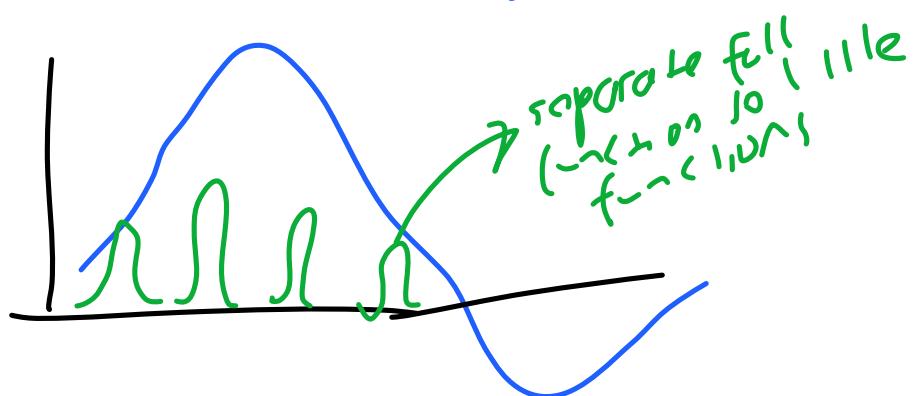
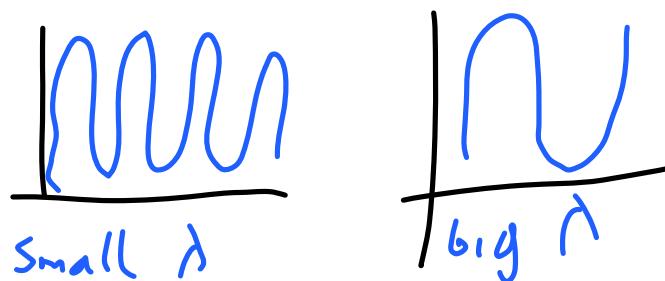
Kernel Regression and RBF (Radial Basis Functions)

15 Mart 2017 Çarşamba 15:00

$$\phi(x) = [K(x, \mu_1, d), K(x, \mu_d, d)]$$

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

Assume that $\mu_1 = 1$ $\mu_d = L$ $M_d = L$ $\lambda = 1$
If you have small lambda, you can
overfit data (multiple signals)

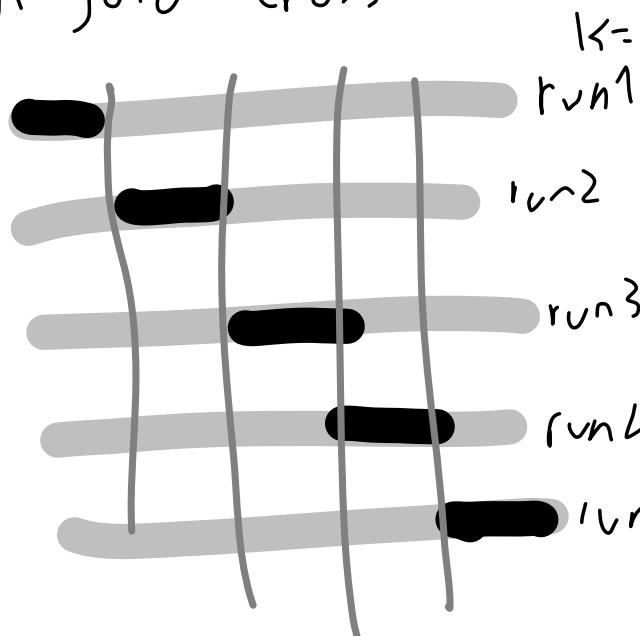


Cross validation

15 Mart 2017 Çarşamba 15:19

Finding best δ

K-fold Cross Validation[^]

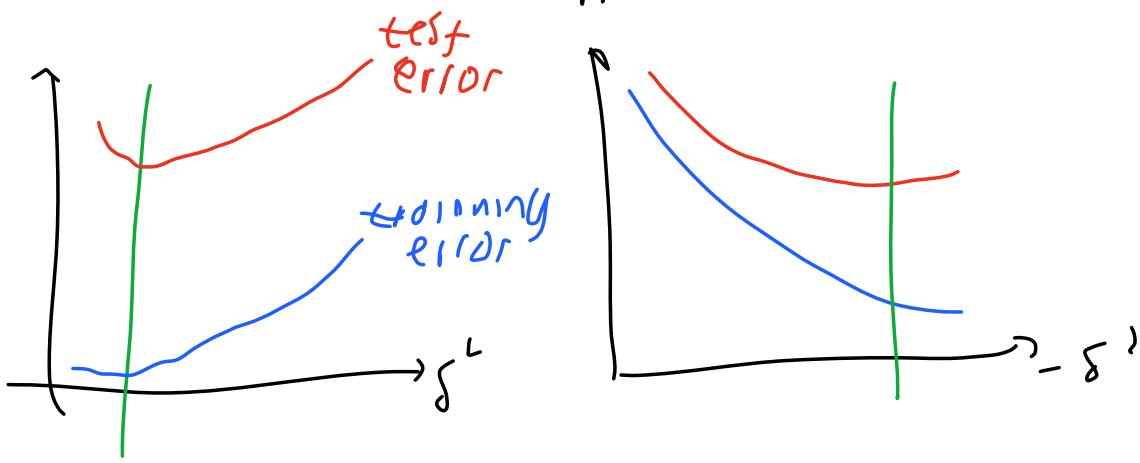


~ 1 the total number
of data samples $k=N$
train TEST run
1 2 3 4 5

train	TEST	run
—	—	1
—	—	2
—	—	3
—	—	4
—	—	5

+ — —
Find for each δ and pick
minimum one

If you have two much classes, this method is hard
and unnecessary to apply



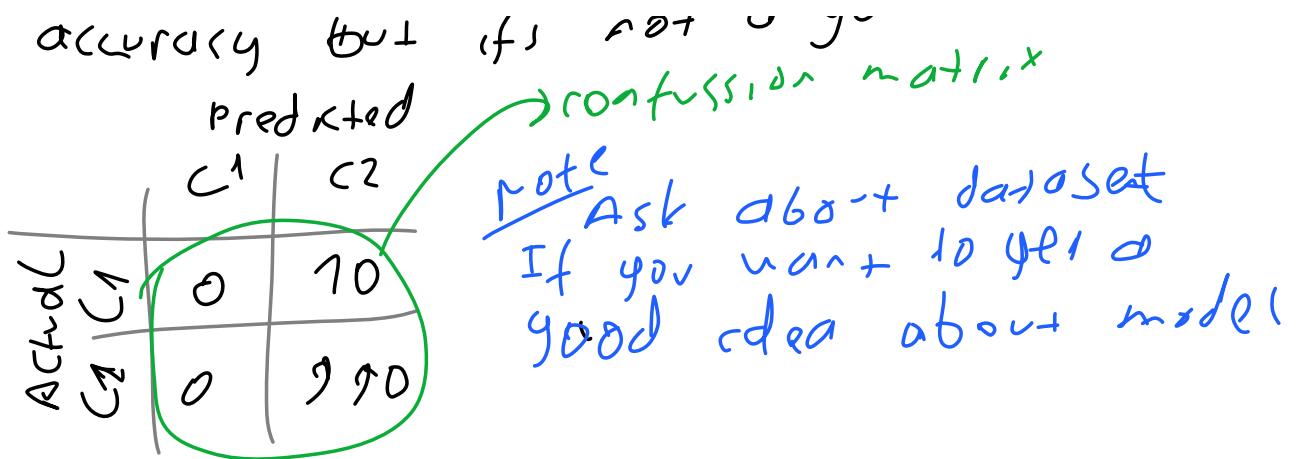
and if

An extreme example

10 C1 990 C2

Unbalanced data set

If you classify everything as C2 you will have 99% accuracy but it's not a good model
and \rightarrow confusion matrix



Bias-Variance Dilemma

22 Mart 2017 Çarşamba 13:42

Generalization Error: Error on the unseen examples. There is a mathematical formula:

Generalization Error: $\text{Bias}^2 + \text{Variance}$

Difference between the actual function and predicted function is named as bias

For example: actual function is linear but you fit 16th order function. Bias will be high. If I will increase the complexity of the model bias will decrease. If we fit 19th order polynomial they wont vary too much from the 16th order polynomial. If we increase complexity of the model, variance will increase.

We don't have a chance to compute bias. Its only possible in simulations.

But we can estimate variance

Since we cant calculate generalization error we will use cross validation

Gradient Descent Optimization Method

22 Mart 2017 Çarşamba 14:09

Used in neural networks.

Closed Form solution: $Q = Q_0 + Q_1 x$ (Q calculated in one step with $(X^T X)^{-1} Y$)

Solution Method

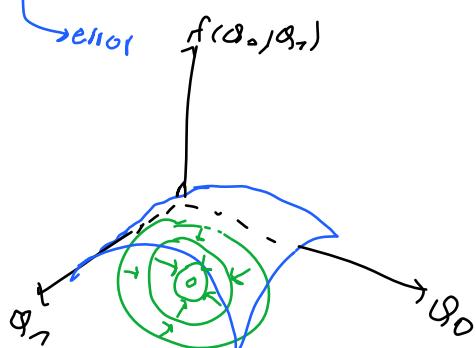
Let Q be a d dimensional vector

$f(Q)$ is a scalar valued function.

Gradient vector of $f(\cdot)$ with respect to $Q = \begin{bmatrix} Q_1 \\ Q_2 \\ \vdots \\ Q_n \end{bmatrix}$

$$\nabla f(Q) = \begin{bmatrix} \frac{\partial f(Q)}{\partial f(Q_1)} \\ \frac{\partial f(Q)}{\partial f(Q_2)} \\ \vdots \\ \frac{\partial f(Q)}{\partial f(Q_d)} \end{bmatrix}$$

$$J(Q) = \sum_{i=1}^n (y_i - x_i Q)^2 \stackrel{d \rightarrow \infty}{\rightarrow} f(Q_0, Q_1)$$



Steepest Gradient Algorithm

-Start with a random Q guess vector

-Update Q *learning rate ex 0.1, 0.2*

$$Q_{k+1} = Q_k -$$

Perception Algorithm

12 Nisan 2017 Çarşamba 13:49

$$x_1^{(1)}, x_2^{(1)}, x_3^{(1)} \dots x_d^{(1)}$$

$$w_0 x_0 + w_1 x_1 + \dots + w_d x_d > thr$$

Always 1 feature

$$w_0 x_0 + w_1 x_1 + \dots + w_d x_d < 0$$

$$w = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_d \end{bmatrix} \quad x = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix}$$

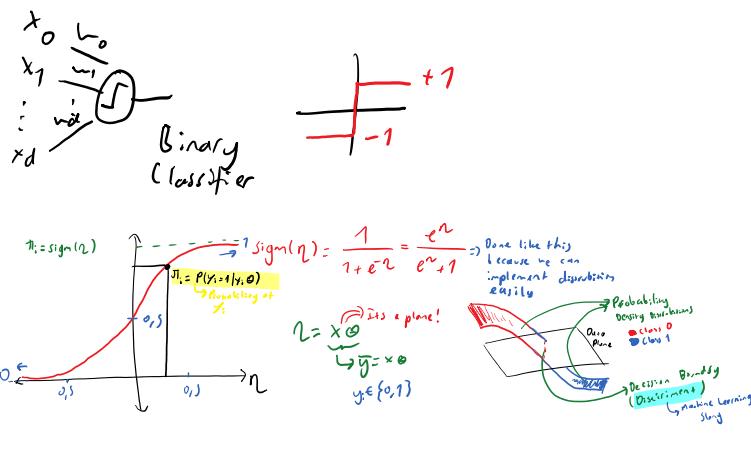
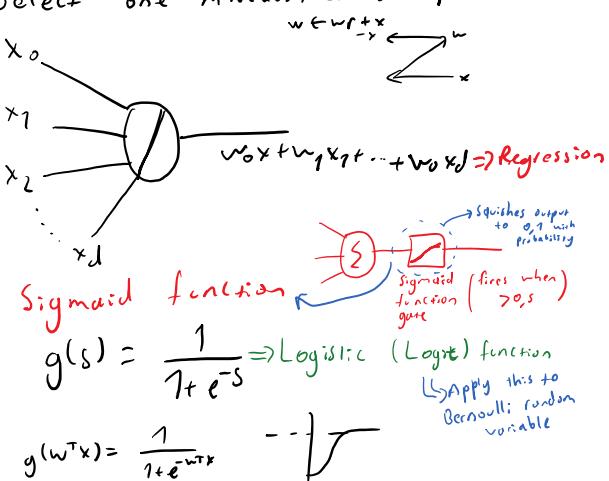
$$r \in \{+1, -1\}$$

if $r > thr$ then +

else -

$$\text{Class label} = \text{sign}(w^T x)$$

Starts with an arbitrary weight vector.
Classify all training examples using $w^T x$
Select one misclassified example



$0, s$ is threshold for classifying

$$y = \frac{1}{1+e^{(w^T x)}} \geq 0, s$$

Using maximum likelihood parameter $\theta = \prod_{i=1}^N (y_i^t)^{r_i^t} (1-y_i^t)^{1-r_i^t} \Rightarrow$ Maximize this for a good model

$$\log \prod_{i=1}^N (y_i^t)^{r_i^t} (1-y_i^t)^{1-r_i^t} = \sum_{i=1}^N \log ((y_i^t)^{r_i^t} (1-y_i^t)^{1-r_i^t}) = \sum_{i=1}^N (r_i^t \log y_i^t + (1-r_i^t) \log (1-y_i^t))$$

this minimize minus instead of maximizing

$-\sum p_i \log p_i \Rightarrow$ Entropy

↳ Information theory \Rightarrow Claude Shannon
John von Neumann
Turing
↳ David Mackay
YouTube (classes)

Bill Gates's book:
Seeing future or such?

$$-\sum r_i^t \log y_i^t + (1-r_i^t) \log (1-y_i^t)$$

Gradient Descent

$$f(\theta) \quad \theta \leftarrow \theta - \eta f(\theta)$$

$$\theta \leftarrow \theta - \eta \nabla f(\theta)$$

$$\nabla(\theta) = \begin{bmatrix} \nabla \theta_1 \\ \nabla \theta_2 \\ \vdots \\ \nabla \theta_l \end{bmatrix}$$

$$\frac{\partial}{\partial \theta_j} \frac{1}{1+e^{w^T x}} = \frac{1}{1+e^{w^T x-w_0}} \quad \text{Start with } w \text{ vector}$$

$$\frac{\partial}{\partial \theta_j} \frac{1}{1+e^{-s}} = \frac{0+e^{-s}}{(1+e^{-s})^2} = \frac{1+e^{-s}-1}{(1+e^{-s})^2} = \frac{1-e^{-s}}{(1+e^{-s})^2} = \frac{1}{1+e^{-s}} - \frac{1}{(1+e^{-s})^2}$$

$$\frac{\partial}{\partial s} g(s) = g(s)(1-g(s))$$

↳ need to add ds (additional multiplication)

$$\frac{\partial}{\partial s} \frac{1}{(1+e^{-s})^2} = \frac{0+e^{-s}}{(1+e^{-s})^2} = \frac{1+e^{-s}-1}{(1+e^{-s})^2} = \frac{1-e^{-s}}{(1+e^{-s})^2} - \frac{1}{(1+e^{-s})^2} = \frac{1}{1+e^{-s}} - \frac{1}{(1+e^{-s})^2} = g(s) (1-g(s))$$

$$\Delta w_j = \sum_{t=1}^N \left[\frac{r^t}{y^t} + \frac{(1-r^t)}{1-y^t} \right] \frac{\partial y^t}{\partial w_j} \Rightarrow \text{using } \frac{\partial}{\partial w} \log s = \frac{1}{s}$$

derivative for each w

$$= \sum_{t=1}^N \left[\frac{r^t}{y^t} \frac{\partial y^t}{\partial w} + \frac{(1-r^t)}{(1-y^t)} \frac{\partial}{\partial w} (1-y^t) \right] = \sum_{t=1}^N \left[\frac{r^t}{y^t} - \frac{(1-r^t)}{1-y^t} \right] \frac{\partial y^t}{\partial w_j}$$

$$\frac{\partial}{\partial w_i} = \frac{1}{1+e^{-(w^T x + w_0)}} = \frac{0+e^{-(w^T x + w_0)}}{1+(e^{-(w^T x + w_0)})^2} \Rightarrow \frac{\partial}{\partial s} g(s) = g(s) (1-g(s)) x_j^T$$

\hookrightarrow outcome of the model

$$w_j \leftarrow w_j - \eta \sum_{t=1}^N \left[\frac{r^t}{y^t} - \frac{(1-r^t)}{1-y^t} \right] y^t (1-y^t) x_j^T$$

$$\Delta w_0 = - \sum_{t=1}^N (r^t - y^t)$$

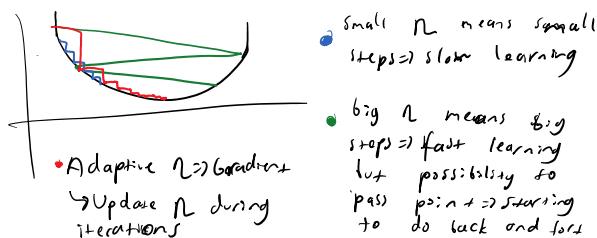
use this two to calculate α values (weights)

At the beginning of the algorithm use very small η 's to stay away from overfitting. Using big random numbers can create phantom dominant or phantom weak features

$r \hat{=} 1$ if $y=1$ 0 if $y=0$ (in the training data)

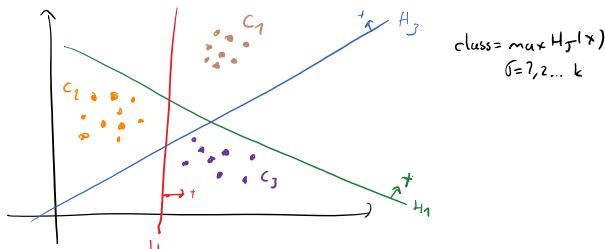
- Find a good η (learning rate)
- Generalize algorithm; after all it's for binary.

An Example Neural Cell:

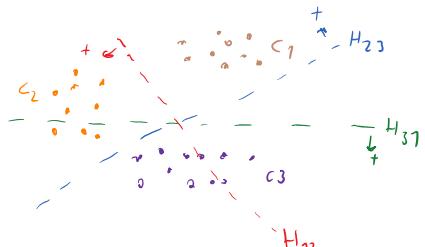


Increasing K

$|K| > 2$?



$$H_{i,2} > 0 \quad c_2 \quad H_{i,1} > 0 \quad c_3 \quad H_{i,2} > 0 \quad c_1$$



choose $c_i : +$

$\forall j \neq i$

$H_{i,j}(v) > 0$

$H_{i,j}(x) > 0$

> 0 if $x \in C_i$

≤ 0 if $x \notin C_i$

sum of

Output of posterior models must be one since:

$$P(C=1 | X) = y^t = \frac{1}{1 + \exp(w^T x + w_0)} \quad r^t \in \{0, 1\}$$

$$P(C=0 | X) = y^0 = \frac{1}{1 + \exp(w^T x + w_0)}$$

$$P(C=0 | X) = 1 - \frac{1}{1 + \exp(w^T x + w_0)}$$

for $k > 2$:

modify function to handle sum of probabilities > 1

$$r^+ \in \text{mult}(k)$$

$$\frac{1}{1+e^{-s}} = \frac{e^s}{1+e^s}$$

$$p(c_i | x) = \frac{\exp(w_i^T x + w_0)}{\sum_{j=1}^k \exp(w_j^T x + w_0)} \quad i=1, \dots, k$$

$$\begin{aligned} p(c=1 | x) \\ p(c=2 | x) \\ \vdots \\ + \frac{p(c=k | x)}{1} \end{aligned}$$

$$(y^+)^{r^+} (1-y^+)^{1-r^+} \Rightarrow r_i^+ = 1 \Rightarrow \text{if class belongs to } c_i \\ \text{we must change those too} \Rightarrow r_i^+ = 0 \Rightarrow \text{else is 0.}$$

$$\begin{aligned} L = \prod_{t=1}^n \prod_{i=1}^k (y_i)^{r_i^+} \\ L = \log \prod_{t=1}^n \prod_{i=1}^k (y_i)^{r_i^+} \Rightarrow \Delta w_{ij} \Rightarrow \text{apply gradient descent} \end{aligned}$$

HW: expand and take derivative
for $w_{i,j}$

Logistic regression: (It's a batch learning \Rightarrow use all data in one time)

$$p(y|x, \theta) = \prod_{i=1}^n \text{Ber}(y_i | \text{sigm}(x_i; \theta))$$

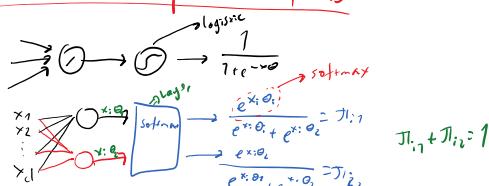
$$= \prod_{i=1}^n \left[\frac{1}{1+e^{-x_i \cdot \theta}} \right]^{y_i} \left[1 - \frac{1}{1+e^{-x_i \cdot \theta}} \right]^{1-y_i} \text{ where } x_i \cdot \theta = \theta_0 + \sum_{j=1}^d \theta_j x_{ij} = \tilde{y}_i$$

y_i is either 1 or 0 so one of the probabilities will not affect result

$$\begin{aligned} c(\theta) &= -\log P(Y|x; \theta) \\ &= -\sum_{i=1}^n y_i \log \theta_i + (1-y_i) \log (1-\theta_i) \Rightarrow \text{cross entropy} \\ \theta_{k+1} &= \theta_k - H^{-1} g_k \\ g_k &= x^T (\pi_k - y) \\ H_k &= x^T S_k x \\ S_k &= \text{diag}(\pi_{1k}(1-\pi_{1k}), \dots, \pi_{dk}(1-\pi_{dk})) \\ \pi_{ik} &= \text{sigm}(x_i \cdot \theta_k) \\ \theta_{k+1} &= (x^T S_k x)^{-1} x^T [S_k x \cdot \theta_k + y - \pi_k] \end{aligned}$$

(next iteration)

Neural Network representation of loss



Indicator

$$\Pi_c(y_i) = \begin{cases} 1 & \text{if } y_i = c \\ 0 & \text{otherwise} \end{cases}$$

$$P(Y|x, \theta) = \prod_{i=1}^n p(y_i | x_i; \theta) = \prod_{i=1}^n \Pi_{i1}(y_i) \Pi_{i2}(y_i)$$

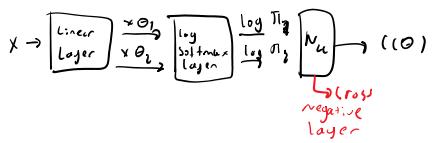
One of them will be ineffective since $\Pi_{i1}=1 \Rightarrow \Pi_{i2}=0$
 $\Pi_{i2}=0 \Rightarrow \Pi_{i1}=1$

$$c(\theta) = -\log P(Y|x, \theta) = -\sum_{i=1}^n \Pi_{i1}(y_i) \log \Pi_{i1} + \Pi_{i2}(y_i) \log \Pi_{i2}$$

Log will make things easier if you are trying to minimize.

$$x \rightarrow \left[\begin{array}{c|c} \text{Linear} & \times \theta \\ \hline \text{Layer} & x \cdot \theta \end{array} \right] \xrightarrow{\text{softmax}} \left[\begin{array}{c|c} \log \Pi_{i1} & \log \Pi_{i2} \\ \hline \text{Layer} & \text{Layer} \end{array} \right] \xrightarrow{\text{Log}} (c(\theta))$$

to minimize.



Modifying Learning Mode with new coming data
(Online Learning)

Learning problem can vary in change. If you apply online learning, learning will be better. Applying Online learning to gradient descent is called Stochastic Gradient Descent.

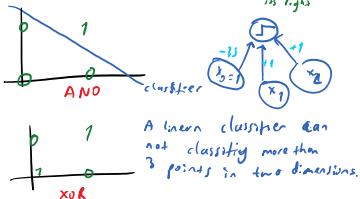
Assume applying to one data:

$$\Delta w = \eta \cancel{\frac{1}{n} (r^t - y^t) x^t} = \eta (r^t - y^t) x^t$$

disappears (one input)
(n=1)

→ observed
→ desired
 $(r^t - y^t) = 0$
in right

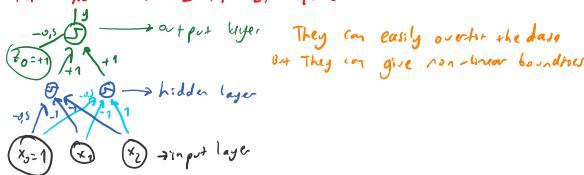
x_1	x_2	AND	xOR
0	0	0	0
0	1	0	1
1	0	0	1
1	1	1	0



For linear classifier in d dimension $d+1$ data points can be correctly classified

Multilayer Neural Network:

For XOR: $x_1 \oplus x_2 = (x_1 \wedge x_2) \vee (\bar{x}_1 \wedge \bar{x}_2)$

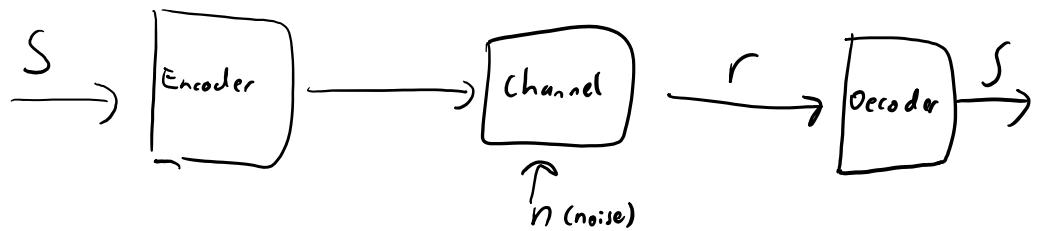


Information Theory

13 Nisan 2017 Perşembe 22:25

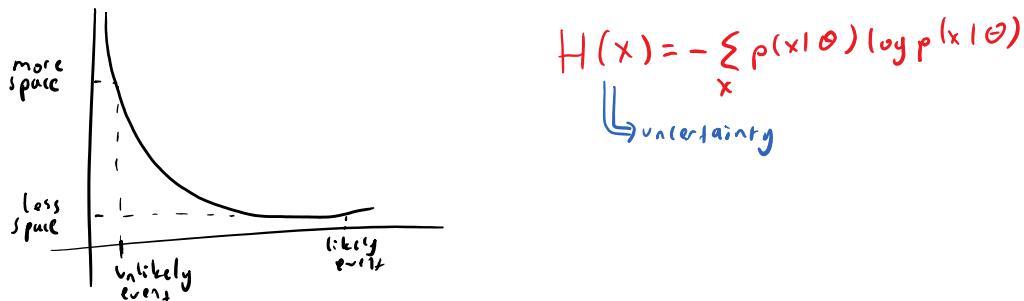
Making a reliable communications systems

- ↳ Change channel to less noisy one
- ↳ Accept noise and handle it



Binary Symmetric Channel (Binary)

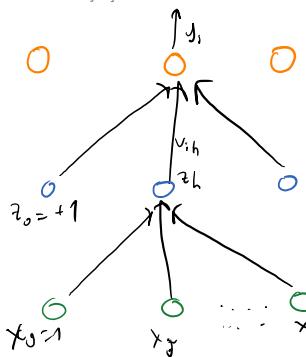
$$\begin{array}{ccc} 0 & \xrightarrow{\quad} & 0 \\ & \cancel{\xrightarrow{\quad}} & \\ 1 & \xrightarrow{\quad} & 1 \end{array} \quad \begin{aligned} p(y=0|x=0) &= 1-f \\ p(y=1|x=0) &= f \end{aligned} \quad \left. \begin{array}{c} \\ \end{array} \right\} \text{Binary Symmetric}$$



Multilayer Perceptrons - chapter 11

19 Nisan 2017 Çarşamba

14:48



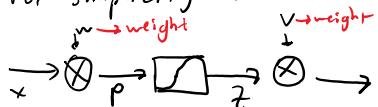
$$y_i = v^T z = \sum_{n=1}^H v_{in} z_n + b_i$$

$$z_h = \text{sigmoid}(w_h^T x) = \frac{1}{1 + \exp[-(\sum_{j=1}^J w_{hj} x_j + b_h)]}$$

An example model for parameter learning will be given and when we will apply backpropagation.

Back propagation

For simplicity we will consider one input & neuron: (can be duplicated)



Using Regression
minimize error ($\epsilon = \frac{1}{2}(r-y)^2$)

Another regression method can be used
If you want to use gradient descent just apply it on L

$$\frac{\partial \epsilon}{\partial v} = \frac{\partial \epsilon}{\partial y} \cdot \frac{\partial y}{\partial v} = -1(r-y) \cdot z \quad \Delta v = -\eta \frac{\partial \epsilon}{\partial v} \text{ update for } v$$

$$\Delta v = \eta (r-y)z \quad v_{i+1} = v_i + \Delta v$$

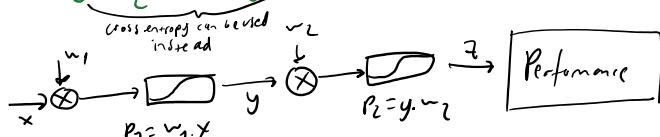
$$\frac{\partial \epsilon}{\partial w} = \frac{\partial \epsilon}{\partial y} \cdot \frac{\partial y}{\partial w} = \frac{\partial \epsilon}{\partial y} \cdot \frac{\partial y}{\partial z} \cdot \frac{\partial z}{\partial w} = \frac{\partial \epsilon}{\partial y} \cdot \frac{\partial y}{\partial z} \cdot \frac{\partial z}{\partial p} \cdot \frac{\partial p}{\partial w}$$

$$-(r-y) v \cdot z \cdot (1-z) \cdot x \quad \Delta w = -\eta \frac{\partial \epsilon}{\partial w}$$

$$\Delta w = \eta (r-y) v \cdot z \cdot (1-z) \cdot x$$

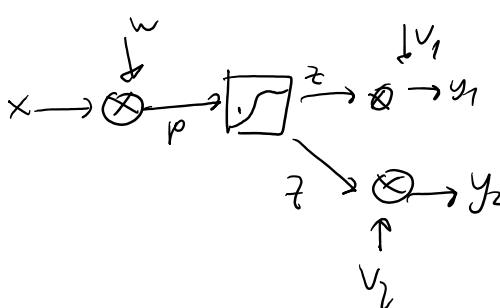
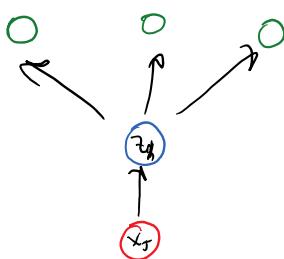
For multiple inputs:

$$\epsilon = \frac{1}{2} \sum (r_i - y_i)^2 \quad \Delta v = \eta \sum \epsilon (r_i - y_i) \cdot z_i$$



$$\log y^r (1-y)^{1-r} \quad \frac{\partial \epsilon}{\partial w_2} = -(r-y) \cdot z \cdot (1-z) \cdot w_2 \cdot y \cdot (1-y) \cdot x$$

Another Model



$$\epsilon = \frac{1}{2} [(r_1 - y_1)^2 + (r_2 - y_2)^2]$$

$$\frac{\partial \epsilon}{\partial v_1} = \frac{\partial \epsilon}{\partial y_1} \cdot \frac{\partial y_1}{\partial v_1} = -(r_1 - y_1) \cdot z$$

$$\frac{\partial \epsilon}{\partial v_2} = - (r_2 - y_2) \cdot z$$

$$\frac{\partial v_1}{\partial y_1} - \frac{\partial v_1}{\partial y_1} \cdot \frac{\partial v_1}{\partial v_1} = -v_1 \cdot v_1 - v_2 \cdot v_2$$

$$= \frac{\partial}{\partial z} \left[\frac{1}{2} (r_1 - y_1)^2 \right] + \frac{\partial}{\partial z} \left[\frac{1}{2} (r_2 - y_2)^2 \right] \quad \frac{\partial \zeta}{\partial z} = \frac{\partial}{\partial z} \left[\frac{1}{2} (r_1 - y_1)^2 + (r_2 - y_2)^2 \right]$$

$$\frac{\partial \zeta}{\partial w} = \frac{\partial \zeta}{\partial p} \cdot \frac{\partial p}{\partial w} = \frac{\partial \zeta}{\partial z} \cdot \frac{\partial z}{\partial p} \cdot \frac{\partial p}{\partial w}$$

(1 - z) x
z(z - 1) x
→ -(r_1 - y_1) v_1 - (r_2 - y_2) v_2

$$\frac{\partial \zeta_1}{\partial z} + \frac{\partial \zeta_2}{\partial z} = \frac{\partial \zeta_1}{\partial y_1} \cdot \frac{\partial y_1}{\partial z} + \frac{\partial \zeta_2}{\partial y_2} \cdot \frac{\partial y_2}{\partial z} = -(r_1 - y_1) \cdot v_1 + -(r_2 - y_2) \cdot v_2$$

$$= [-(r_1 - y_1) v_1 - (r_2 - y_2) v_2]$$

$$= -\sum_k (r_k - y_k) v_k \quad \text{for } k \text{ output}$$



Lecture Slides for
INTRODUCTION TO
Machine Learning
2nd Edition
CHAPTER 11:
Multilayer Perceptrons
ETHEM ALPAYDIN
© The MIT Press, 2010
Edited and expanded for CS 4641 by Chris Simpkins
alpaydin@boun.edu.tr
<http://www.cmpe.boun.edu.tr/~ethem/l2ml2e>

Overview

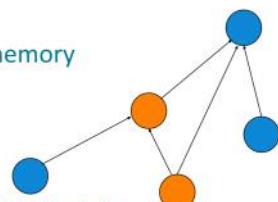
- Neural networks, brains, and computers
- Perceptrons
 - Training
 - Classification and regression
 - Linear separability
- Multilayer perceptrons
 - Universal approximation
 - Backpropagation

Lecture Notes for E Alpaydin 2004 Introduction
to Machine Learning © The MIT Press (V1.1)

2

Neural Networks

- Networks of processing units (neurons) with connections (synapses) between them
- Large number of neurons: 10^{10}
- Large connectivity: 10^5
- Parallel processing
- Distributed computation/memory
- Robust to noise, failures



Lecture Notes for E Alpaydin 2004 Introduction
to Machine Learning © The MIT Press (V1.1)

3

Understanding the Brain

- Levels of analysis (Marr, 1982)
 1. Computational theory
 2. Representation and algorithm
 3. Hardware implementation
- Reverse engineering: From hardware to theory
- Parallel processing: SIMD vs MIMD
Neural net: SIMD with modifiable local memory
Learning: Update by training/experience

Lecture Notes for E Alpaydin 2004 Introduction
to Machine Learning © The MIT Press (V1.1)

4

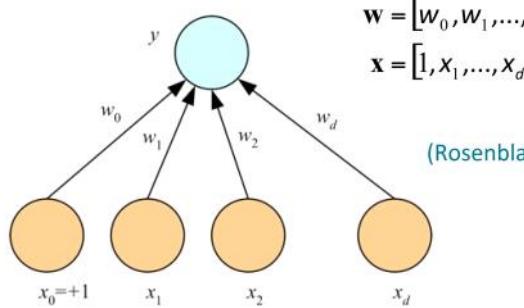
Perceptron

$$y = \sum_{j=1}^d w_j x_j + w_0 = \mathbf{w}^T \mathbf{x}$$

$$\mathbf{w} = [w_0, w_1, \dots, w_d]$$

$$\mathbf{x} = [1, x_1, \dots, x_d]$$

(Rosenblatt, 1962)



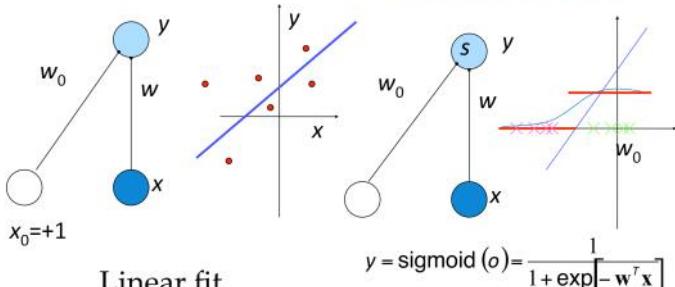
Lecture Notes for E Alpaydin 2004 Introduction
to Machine Learning © The MIT Press (V1.1)

5

What a Perceptron Does

- Regression: $y = \mathbf{w}^T \mathbf{x} + w_0$

- Classification: $y = 1(\mathbf{w}^T \mathbf{x} + w_0 > 0)$



Lecture Notes for E Alpaydin 2004 Introduction
to Machine Learning © The MIT Press (V1.1)

6

K Outputs

Regression:

$$y_i = \sum_{j=1}^d w_{ij} x_j + w_{i0} = \mathbf{w}_i^T \mathbf{x}$$

$$\mathbf{y} = \mathbf{Wx}$$

Classification:

$$o_i = \mathbf{w}_i^T \mathbf{x}$$

$$y_i = \frac{\exp o_i}{\sum_k \exp o_k}$$

choose C_i
if $y_i = \max_k y_k$

Lecture Notes for E Alpaydin 2004 Introduction to Machine Learning © The MIT Press (V1.1)

Training

- Online (instances seen one by one) vs batch (whole sample) learning:
 - No need to store the whole sample
 - Problem may change in time
 - Wear and degradation in system components
- Stochastic gradient-descent: Update after a single pattern
- Generic update rule (LMS rule):

$$\Delta w_{ij}^t = \eta (r_i^t - y_i^t) x_j$$

Update=LearningFactor * DesiredOutput-AcualOutput * Input

Lecture Notes for E Alpaydin 2004 Introduction to Machine Learning © The MIT Press (V1.1)

Training a Perceptron: Regression

- Regression (Linear output):

$$E^t(\mathbf{w} | \mathbf{x}^t, r^t) = \frac{1}{2} (r^t - y^t)^2 = \frac{1}{2} [r^t - (\mathbf{w}^T \mathbf{x}^t)]^2$$

$$\Delta w_j^t = \eta (r^t - y^t) x_j^t$$

Lecture Notes for E Alpaydin 2004 Introduction to Machine Learning © The MIT Press (V1.1)

Classification

- Single sigmoid output

$$y^t = \text{sigmoid}(\mathbf{w}^T \mathbf{x}^t)$$

$$E^t(\mathbf{w} | \mathbf{x}^t, \mathbf{r}^t) = -r^t \log y^t - (1 - r^t) \log (1 - y^t)$$

$$\Delta w_j^t = \eta (r^t - y^t) x_j^t$$

- $K > 2$ softmax outputs

$$y^t_i = \frac{\exp \mathbf{w}_i^T \mathbf{x}^t}{\sum_k \exp \mathbf{w}_k^T \mathbf{x}^t} \quad E^t(\{\mathbf{w}_i\} | \mathbf{x}^t, \mathbf{r}^t) = -\sum_i r_i^t \log y_i^t$$

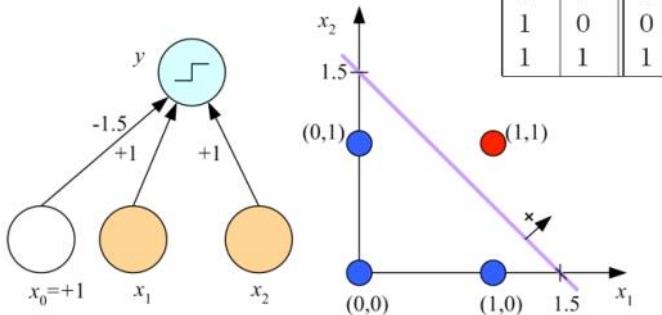
$$\Delta w_{ij}^t = \eta (r_i^t - y_i^t) x_j^t$$

Same as for linear discriminants from chapter 10 except we update after each instance

Lecture Notes for E Alpaydin 2004 Introduction to Machine Learning © The MIT Press (V1.1)

10

Learning Boolean AND

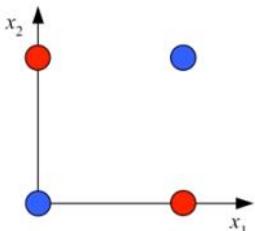


Lecture Notes for E Alpaydin 2004 Introduction to Machine Learning © The MIT Press (V1.1)

11

XOR

x_1	x_2	r
0	0	0
0	1	1
1	0	1
1	1	0



- No w_0, w_1, w_2 satisfy:

$$w_0 \leq 0 \quad (\text{Minsky and Papert, 1969})$$

$$w_2 + w_0 > 0$$

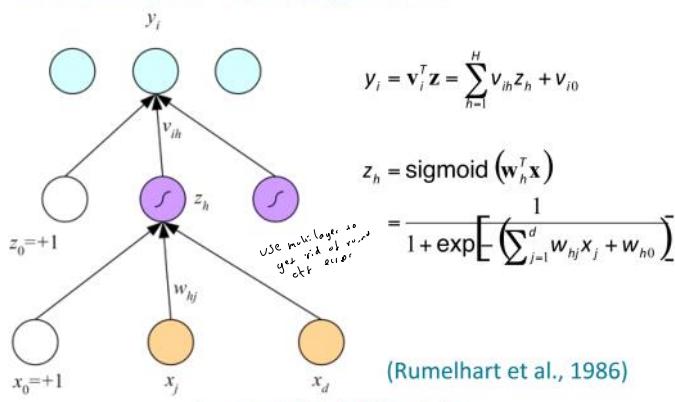
$$w_1 + w_0 > 0$$

$$w_1 + w_2 + w_0 \leq 0$$

Lecture Notes for E Alpaydin 2004 Introduction to Machine Learning © The MIT Press
(V1.1)

12

Multilayer Perceptrons

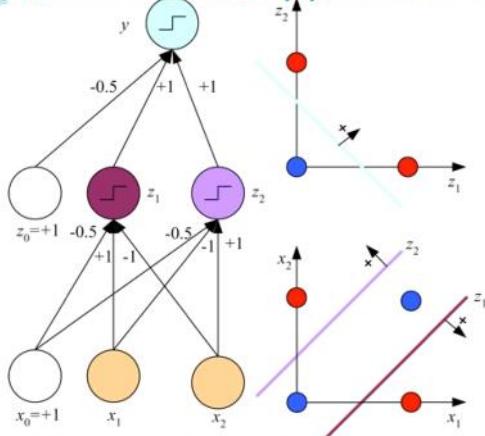


Lecture Notes for E Alpaydin 2004 Introduction to Machine Learning © The MIT Press (V1.1)

important

13

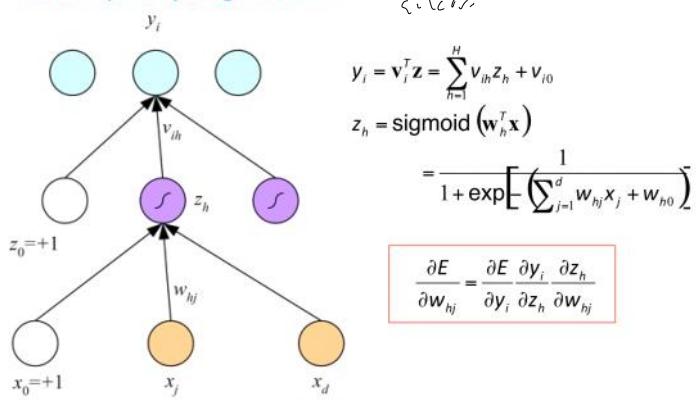
MLP as Universal Approximator



Lecture Notes for E Alpaydin 2004 Introduction to Machine Learning © The MIT Press (V1.1)

14

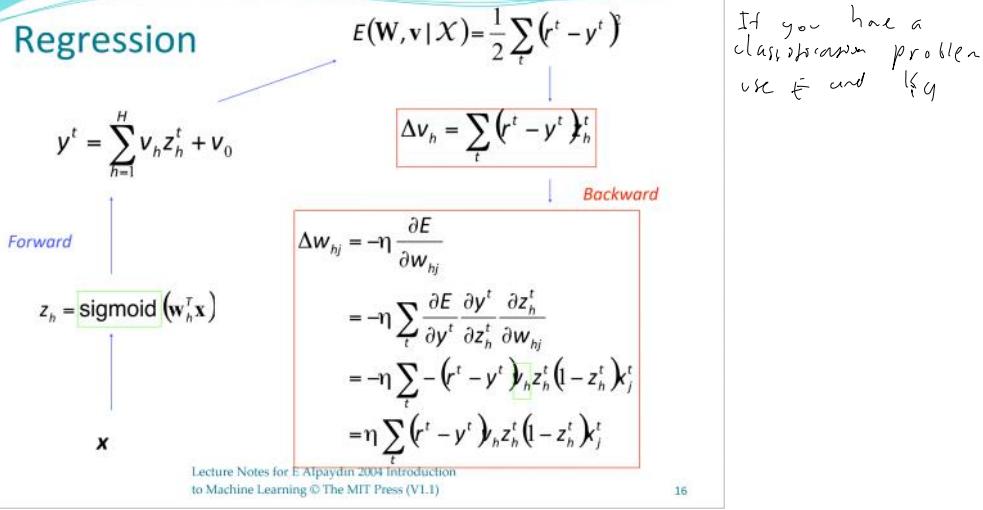
Backpropagation



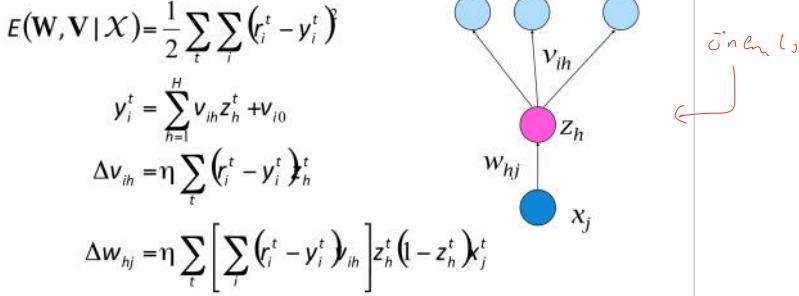
Lecture Notes for E Alpaydin 2004 Introduction to Machine Learning © The MIT Press (V1.1)

15

Regression



Regression with Multiple Outputs



```

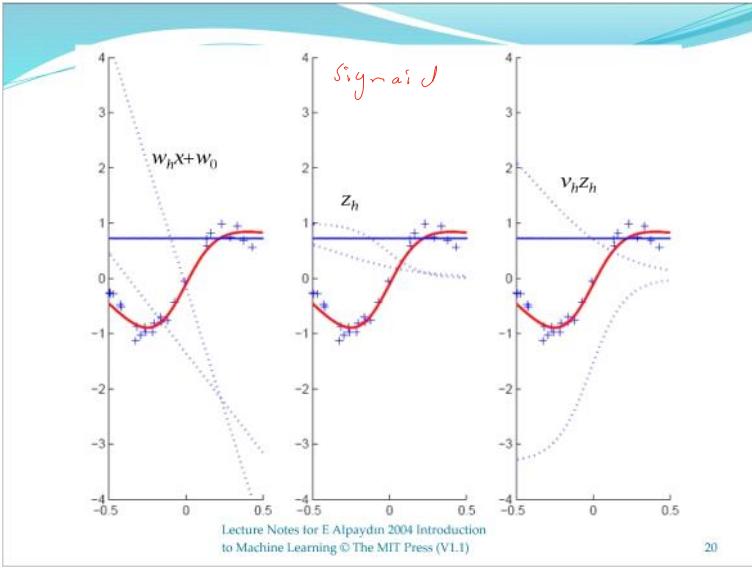
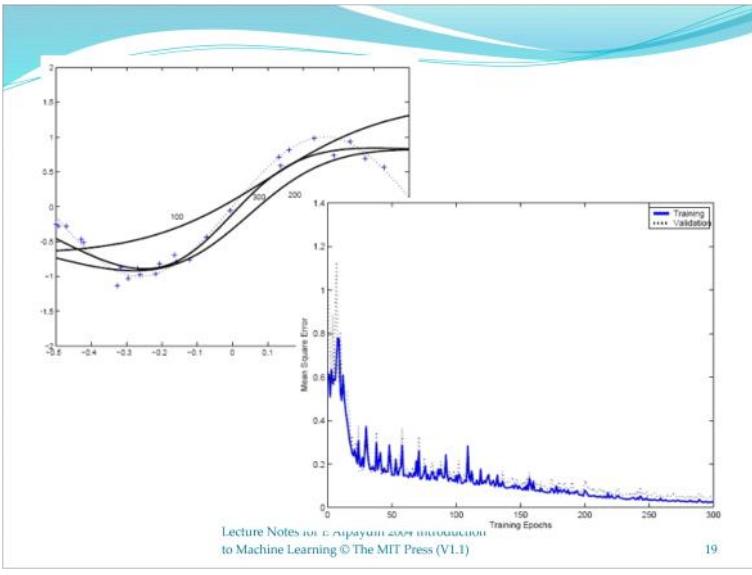
Initialize all  $v_{ih}$  and  $w_{hj}$  to  $\text{rand}(-0.01, 0.01)$ 
Repeat
  For all  $(\mathbf{x}^t, r^t) \in \mathcal{X}$  in random order
    For  $h = 1, \dots, H$ 
       $z_h \leftarrow \text{sigmoid}(\mathbf{w}_h^T \mathbf{x}^t)$ 
    For  $i = 1, \dots, K$ 
       $y_i = \mathbf{v}_i^T \mathbf{z}$ 
    For  $i = 1, \dots, K$ 
       $\Delta \mathbf{v}_i = \eta(r_i^t - y_i^t) \mathbf{z}$ 
    For  $h = 1, \dots, H$ 
       $\Delta \mathbf{w}_h = \eta \left( \sum_i (r_i^t - y_i^t) v_{ih} \right) z_h (1 - z_h) \mathbf{x}^t$ 
    For  $i = 1, \dots, K$ 
       $\mathbf{v}_i \leftarrow \mathbf{v}_i + \Delta \mathbf{v}_i$ 
    For  $h = 1, \dots, H$ 
       $\mathbf{w}_h \leftarrow \mathbf{w}_h + \Delta \mathbf{w}_h$ 
Until convergence  $\Rightarrow$  error rate

```

Apply $\Delta \mathbf{v}_i$ in \mathbf{v}_i

to Machine Learning © The MIT Press (V1.1)

18



Two-Class Discrimination

- One sigmoid output y^t for $P(C_1 | \mathbf{x}^t)$ and $P(C_2 | \mathbf{x}^t) \equiv 1 - y^t$

$$y^t = \text{sigmoid} \left(\sum_{h=1}^H v_h z_h^t + v_0 \right)$$

$$E(W, v | X) = - \sum_t r^t \log y^t + (1 - r^t) \log (1 - y^t)$$

$$\Delta v_h = \eta \sum_t (r^t - y^t) z_h^t$$

$$\Delta w_{hj} = \eta \sum_t (r^t - y^t) z_h^t (1 - z_h^t) x_j^t$$

K>2 Classes

$$o_i^t = \sum_{h=1}^H v_{ih} z_h^t + v_{i0} \quad y_i^t = \frac{\exp o_i^t}{\sum_k \exp o_k^t} = P(C_i | \mathbf{x}^t)$$
$$E(W, v | X) = - \sum_t \sum_i r_i^t \log y_i^t$$
$$\Delta v_{ih} = \eta \sum_t (r_i^t - y_i^t) z_h^t$$
$$\Delta w_{hj} = \eta \sum_t \left[\sum_i (r_i^t - y_i^t) y_{ih} \right] z_h^t (-z_h^t) z_j^t$$

Lecture Notes for E Alpaydin 2004 Introduction to Machine Learning © The MIT Press (V1.1)

22

Layer [tiplication]

because of the numerical issues every layer learned individually and then combined

Multiple Hidden Layers

- MLP with one hidden layer is a **universal approximator** (Hornik et al., 1989), but using multiple layers may lead to simpler networks

$$z_{1h} = \text{sigmoid}(\mathbf{w}_{1h}^T \mathbf{x}) = \text{sigmoid}\left(\sum_{j=1}^d w_{1hj} x_j + w_{1h0}\right) h = 1, \dots, H_1$$

$$z_{2l} = \text{sigmoid}(\mathbf{w}_{2l}^T \mathbf{z}_1) = \text{sigmoid}\left(\sum_{h=1}^{H_1} w_{2lh} z_{1h} + w_{2l0}\right) l = 1, \dots, H_2$$

$$y = \mathbf{v}^T \mathbf{z}_2 = \sum_{l=1}^{H_2} v_l z_{2l} + v_0$$

Lecture Notes for E Alpaydin 2004 Introduction to Machine Learning © The MIT Press (V1.1)

23

Improving Convergence

- Momentum

$$\Delta w_i^t = -\eta \frac{\partial E^t}{\partial w_i} + \alpha \Delta w_i^{t-1}$$

- Adaptive learning rate

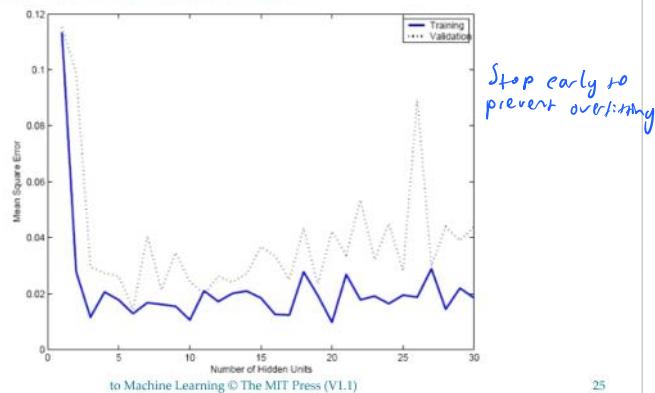
$$\Delta \eta = \begin{cases} +a & \text{if } E^{t+\tau} < E^t \\ -b\eta & \text{otherwise} \end{cases}$$

Lecture Notes for E Alpaydin 2004 Introduction to Machine Learning © The MIT Press (V1.1)

24

Overfitting/Overtraining

Number of weights: $H(d+1)+(H+1)K$

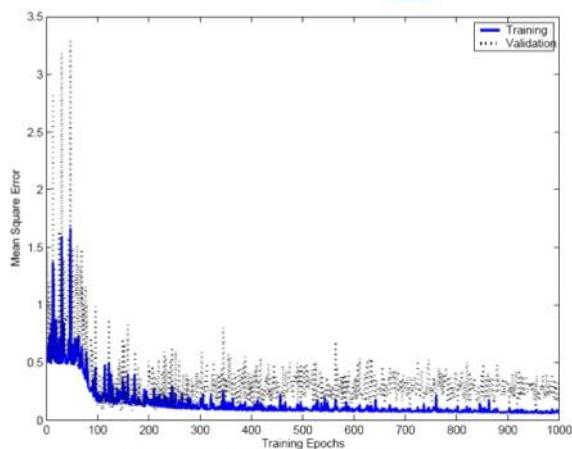


Conclusion

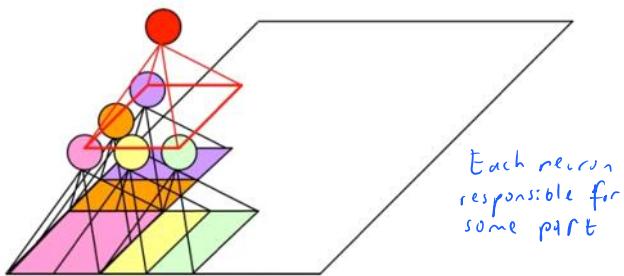
- Perceptrons handle linearly separable problems
- Multilayer perceptrons handle any problem
- Logistic discrimination functions enable gradient descent-based backpropagation
 - Solves the structural credit assignment problem
 - Susceptible to local optima
 - Susceptible to overfitting

Lecture Notes for E Alpaydin 2004 Introduction to Machine Learning © The MIT Press (V1.1)

26



Structured MLP

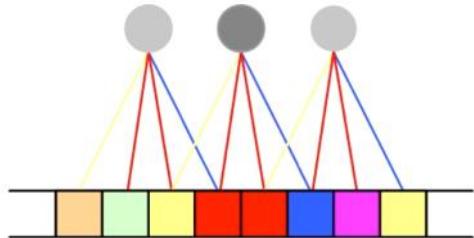


(Le Cun et al, 1989)

Lecture Notes for E Alpaydin 2004 Introduction
to Machine Learning © The MIT Press (V1.1)

28

Weight Sharing



Lecture Notes for E Alpaydin 2004 Introduction
to Machine Learning © The MIT Press (V1.1)

29

Hints

(Abu-Mostafa, 1995)

- Invariance to translation, rotation, size



- Virtual examples

- Augmented error: $E' = E + \lambda_h E_h$

If x' and x are the "same": $E_h = [g(x|\theta) - g(x'|\theta)]^2$

Approximation hint:

$$E_h = \begin{cases} 0 & \text{if } g(x|\theta) \in [a_x, b_x] \\ (g(x|\theta) - a_x)^2 & \text{if } g(x|\theta) < a_x \\ (g(x|\theta) - b_x)^2 & \text{if } g(x|\theta) > b_x \end{cases}$$

Lecture Notes for E Alpaydin 2004 Introduction
to Machine Learning © The MIT Press (V1.1)

30

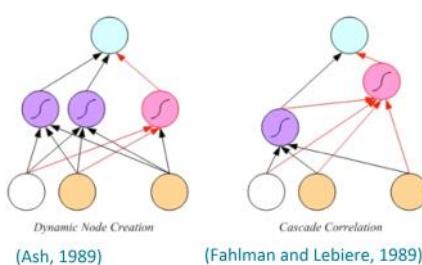
Tuning the Network Size

→ If weights are small destroy it

- Destructive
- Weight decay:
- Constructive
- Growing networks

$$\Delta w_i = -\eta \frac{\partial E}{\partial w_i} - \lambda w_i$$

$$E' = E + \frac{\lambda}{2} \sum_i w_i^2$$



Lecture Notes for E Alpaydin 2004 Introduction to Machine Learning © The MIT Press (V1.1)

31

Bayesian Learning

- Consider weights w_i as random vars, prior $p(w_i)$
- $$p(\mathbf{w} | X) = \frac{p(X | \mathbf{w}) p(\mathbf{w})}{p(X)} \quad \hat{\mathbf{w}}_{MAP} = \arg \max_{\mathbf{w}} \log p(\mathbf{w} | X)$$
- $$\log p(\mathbf{w} | X) = \log p(X | \mathbf{w}) + \log p(\mathbf{w}) + C$$
- $$p(\mathbf{w}) = \prod_i p(w_i) \text{ where } p(w_i) = c \exp \left[-\frac{w_i^2}{2(1/2\lambda)} \right]$$
- $$E' = E + \lambda \|\mathbf{w}\|^2$$

- Weight decay, ridge regression, regularization
cost = data-misfit + λ complexity

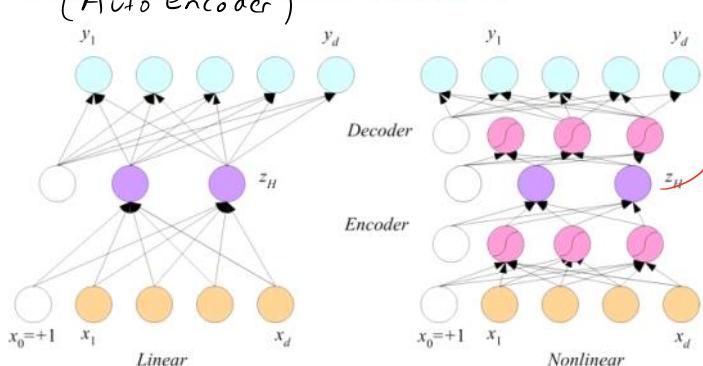
More about Bayesian methods in chapter 14

Lecture Notes for E Alpaydin 2004 Introduction to Machine Learning © The MIT Press (V1.1)

32

Dimensionality Reduction

(Auto Encoder)



Lecture Notes for E Alpaydin 2004 Introduction to Machine Learning © The MIT Press (V1.1)

Learning from decoder and after, using middle layer you can feed classifiers

PCA — Uses SVD

Principle Component Analysis
dimensional reduction

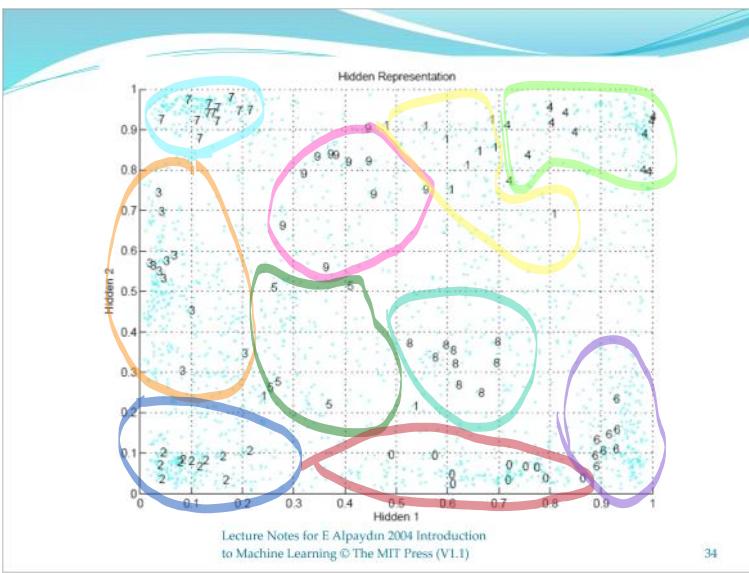
to avoid curse of dimensionality
Rule of thumb (in 2D) used. This also might reduce noise.

process is to 1 dimension
where variance maximized
also used in visualization

Maximize
 $(\mathbf{W}^T \mathbf{x} - \mathbf{E}[\mathbf{w}^T \mathbf{x}])$
 \mathbf{X} → Training data set

Subtract mean from each instance
to increase a zero mean data
Then compute $\mathbf{x}^T \mathbf{x}$ (covariance matrix $(d \times d)$)
compute the eigen vectors of cov. matrix
We have d eigen values and d eigen vectors
Largest of the eigen values will be dominant
Select k biggest α

eigenvalues/eigen vectors
dimensions desired



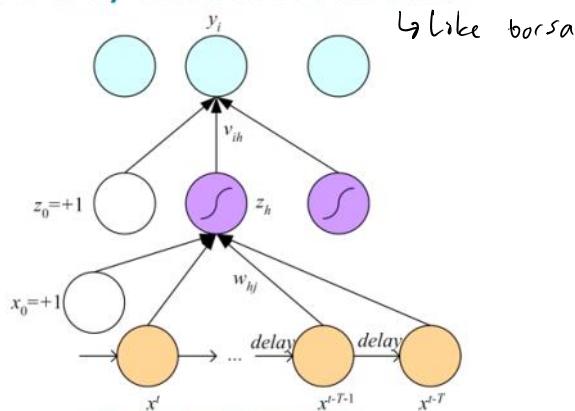
Learning Time

- Applications:
 - Sequence recognition: Speech recognition
 - Sequence reproduction: Time-series prediction
 - Sequence association
- Network architectures
 - Time-delay networks (Waibel et al., 1989)
 - Recurrent networks (Rumelhart et al., 1986)

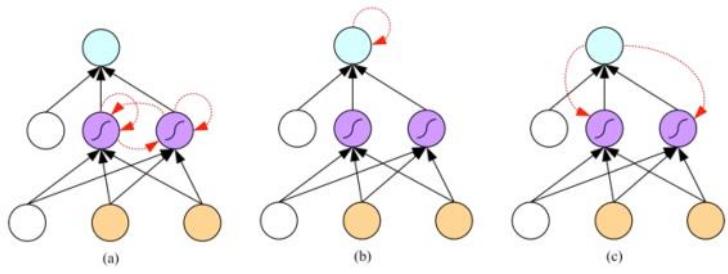
Lecture Notes for E Alpaydin 2004 Introduction to Machine Learning © The MIT Press (V1.1)

35

Time-Delay Neural Networks



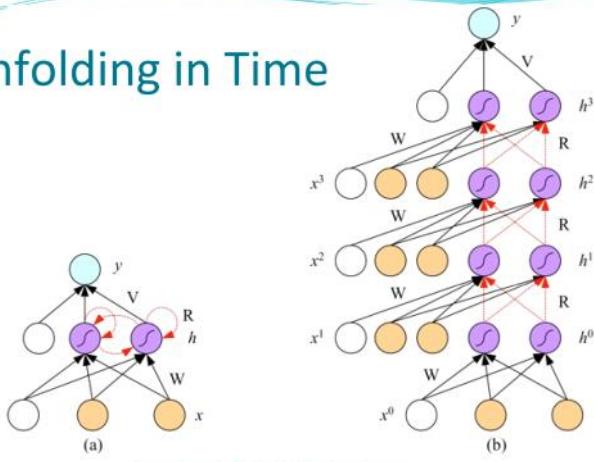
Recurrent Networks



Lecture Notes for E Alpaydin 2004 Introduction
to Machine Learning © The MIT Press (V1.1)

37

Unfolding in Time

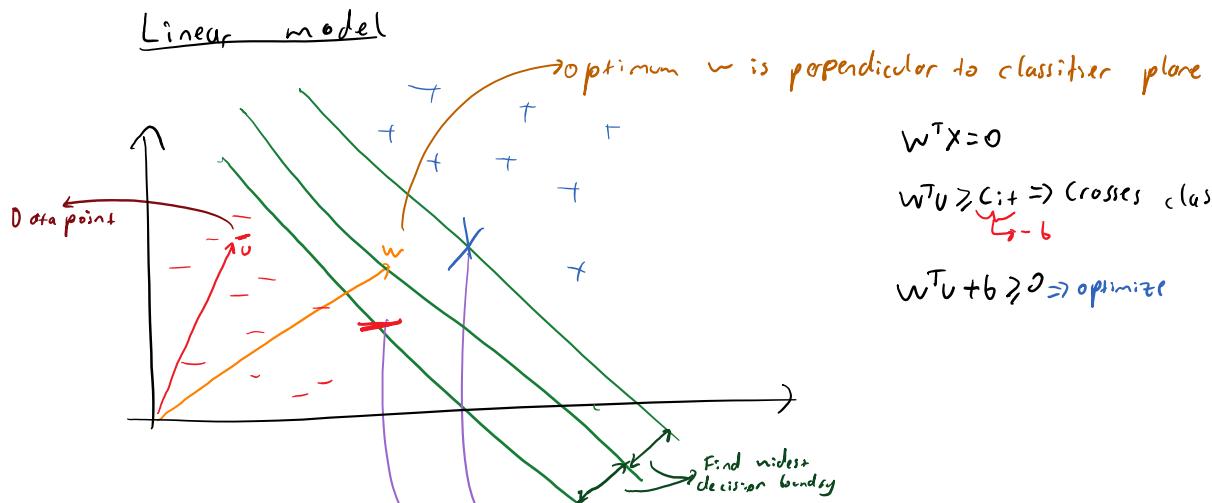


Lecture Notes for E Alpaydin 2004 Introduction
to Machine Learning © The MIT Press (V1.1)

38

Support Vector Machines

3 Mayıs 2017 Çarşamba 13:48



$$w^T x = 0$$

$w^T u \geq c_i + \epsilon$ ⇒ crosses classifier line. It's a plus

$w^T u + b \geq 0$ ⇒ optimize

Support
Vectors
(will be calculated
and will be used to
calculate w)

$$\begin{aligned} y=1 \Rightarrow w^T x_+ + b \geq 1 \\ y=-1 \Rightarrow w^T x_- + b \leq -1 \end{aligned}$$

on the line $\Rightarrow w^T x_+ + b = 1$
(support vectors) $\Rightarrow w^T x_- + b = -1$

Combine two classifiers

$$\begin{aligned} y_i (w^T x_i + b) \geq 1 \\ y_i (w^T x_i + b) - 1 \geq 0 \end{aligned}$$

↳ Data points on the border
(support vectors)

$$w \cdot d \cdot h(d) = (x_+ - x_-)^T \cdot \frac{w}{\|w\|} = \frac{x_+^T w - x_-^T w}{\|w\|}$$

$$x_+^T w = \frac{(1-b)}{\|w\|} - \frac{(-b-1)}{\|w\|} = \frac{1-b+b+1}{\|w\|} = \frac{2}{\|w\|}$$

$$\text{minimize } \frac{1}{\|w\|} \Rightarrow \text{maximize } \|w\| = \sqrt{w_1^2 + w_2^2}$$

$$\text{minimize } \frac{1}{2} w^T w = (\|w\|_2)^2$$

Since we have constraints (More than one function to satisfy)
Lagrange Optimization will be used.

Lagrange multipliers

Max f(x) with respect to g(x)=b

$$L(x, \lambda) = f(x) - \lambda(g(x) - b) \quad \lambda \geq 0$$

$$g(x) = y_i (w^T x_i + b) - 1$$

$$\mathcal{L} = \frac{1}{2} w^T w - \sum \lambda_i [y_i (w^T x_i + b) - 1] \Rightarrow \text{primal form}$$

↳ Lagrange multiplier for each data point (λ_i)

→ maximize

$$\frac{\partial L}{\partial w} = w - \sum \lambda_i y_i x_i = 0 \quad w = \sum \lambda_i y_i x_i$$

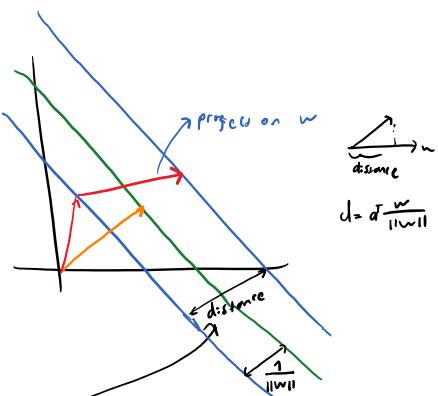
$$L_d = \frac{1}{2} \left(\sum \lambda_i y_i x_i \right)^T \left(\sum \lambda_j y_j x_j \right) - \sum \lambda_i y_i x_i^T \left(\sum \lambda_j y_j x_j \right) - \sum \lambda_i y_i b + \sum \lambda_i$$

$$\frac{\partial L}{\partial b} = \sum \lambda_i y_i = 0 \quad L_d = \sum \lambda_i - \frac{1}{2} \sum \lambda_i \sum \lambda_j y_i y_j x_i^T x_j$$

$$(x_1^T + x_2^T)(x_1 + x_2) = x_1^T x_1 + x_1^T x_2 + x_2^T x_1 + x_2^T x_2$$

T ~ F:nd

Quadratic



Quadratic Optimization tools

$$(x_1^T + x_2^T)(x_1 + x_2) = x_1^T x_1 + x_1^T x_2 + x_2^T x_1 + x_2^T x_2$$

$$\max_{\alpha} L_d = -\frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j x_i^T x_j \Rightarrow \text{find maximum}$$

Decision function: $w^T u + b \geq 0$

$x_i^T x_j = K(x_i, x_j)$

vector to be classified

Previous problem:

$\sum \alpha_i y_i x_i^T u + b \geq 0$

It becomes linearly separable (2 dimension)

$z = \phi x$

$x_i^T u \Rightarrow \phi^T(x_i) \phi(u)$

$x_i^T x_j \Rightarrow \phi^T(x_i) \phi(x_j)$

$\phi(x_i) \phi(x_j) = K(x_i, x_j)$

called "Kernel function"

$K(x, y) = (x^T y + 1)^2 \rightarrow \text{polynomial kernel}$

$\phi(x) = [1, x_1 y_1, x_1^2, x_2^2]$

$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad \phi x \rightarrow \phi \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$

An example ϕ function

ϕ defined according to Mercer's conditions

$K(x, y) = \exp(-\frac{\|x-y\|}{2\sigma^2}) \Rightarrow \text{Radial Basis Function (RBF)}$

Properties of the kernel

Finding α 's
Coordinate Ascents
 $\max_w w(\alpha_1, \alpha_2, \dots, \alpha_m)$

Pseudocode:

Loop until convergence
for $i=1, \dots, m$
 $\alpha_i = \arg \max w(\alpha_1, \dots, \alpha_{i-1}, \alpha_i, \alpha_{i+1}, \dots, \alpha_m)$

Sequential Minimal Optimization \Rightarrow tool box \Rightarrow libSVM

Repeat Until Converge
Select some pair α_i, α_j to update next
Reoptimize $w(\alpha)$ with respect to α_i and α_j
while holding all the other α_k ($k \neq i, j$)

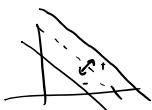
$$\alpha_i + \alpha_j = \sum_{k=1}^m \alpha_k$$

If done right some of the α 's are going to be 0
 $\alpha \neq 0$ are going to be support vectors.

In the book: Kernel Machines

↳ Prof. manyan kimse yok
↳ Andrew Ng machine learning
↳ MIT-AT-SVM

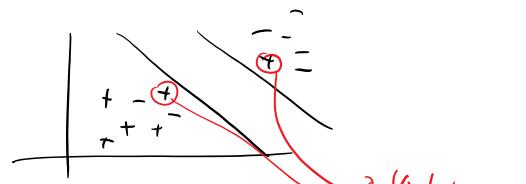
$$\sum \alpha_i y_i x_i^T u + b \geq 0$$



If data is not linearly separable:

$$L_p = \frac{1}{2} w^T w - \sum_i [y_i (w^T x_i + b) - 1] + C \sum_i [y_i - \sum d_i] \rightarrow h(x)$$

$y(x) \rightarrow$ punish violators



Violate boundary
Add constraint

$$g(x) = b$$

$$h(x) = C$$

$g(x) > 0 \Rightarrow$ constraint

$$\begin{aligned}L(x, \lambda) &= f(x) - \lambda(g(x) - b) \\&= p_h(x) - c\end{aligned}$$

Hierarchical Clustering

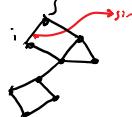
10 Mayıs 2017 Çarşamba 13:46

K-means, Expectation maximization

↳ works well with:
Mixture of gaussians
But to complicated
optimization (?) - check

Spectral Clustering

Considers the similarities between data samples. It uses graph partitioning with constructed & Affinity matrix.



Hierarchical Clustering

Based on similarities and distances

L_p norm L₂ → Euclidean distance

Minkowski distance $d_m(x, y) = \left(\sum_{i=1}^d (x_i - y_i)^p \right)^{1/p}$

City Block (Manhattan) distance = $d_{\text{manhattan}}(x, y) = \sum_{i=1}^d |x_i - y_i|$

Agglomerative clustering

Each data sample assumed to be cluster.

Merge two closest groups at each iteration.

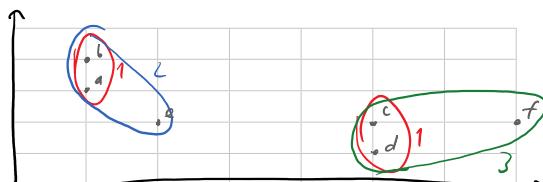
Iterate until there is one group.

Distance between groups:

1) Single linkage $d(G_i, G_j) = \min_{x \in G_i, y \in G_j} (x, y)$

2) Complete link $d(G_i, G_j) = \max_{x \in G_i, y \in G_j} (x, y)$

3) Average-link, centroid

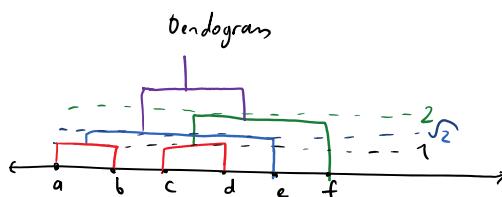


Group Data
- Use euclidean distance
 $a-b: 1$
 $c-d: 1$ } combine

Single linkage
- Use min distance (Consider closest data points to calculate)
 $a-c: \sqrt{2}$

Complete linkage
- Use max distance (Consider most far data points.)

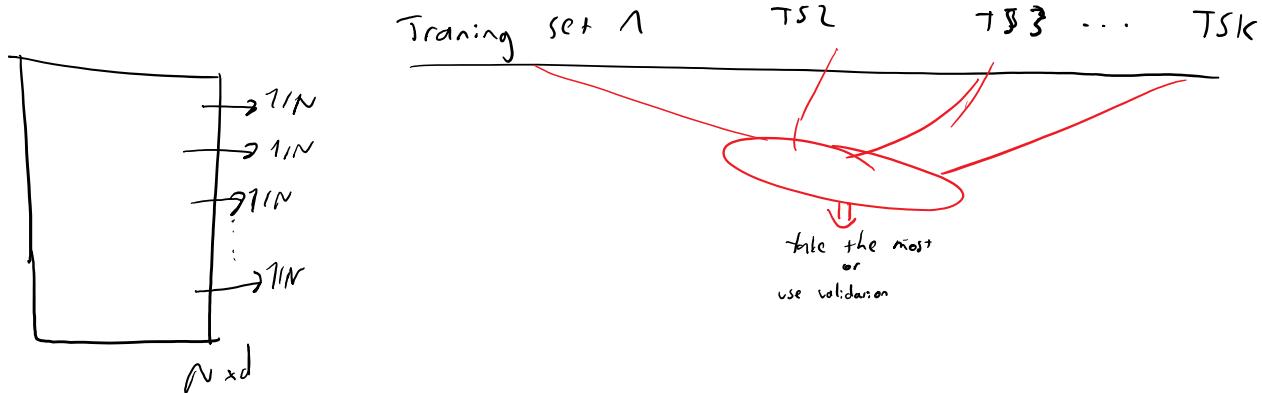
Average Link, centroid
- Use Av. distance (Consider mean of cluster)



Classification Combination (Classifier Ensembles)

10 Mayıs 2017 Çarşamba 15:04

Bagging



Ada Boost

Diversity de incelenir

Unstable algoritmalar daha iyi galisir.