

$K = X \cdot X^T$ linear kernel

$$K_{ij} = K_{ji} = \langle x_j, x_i \rangle = x_j \cdot x_i^T = x_i \cdot x_j^T$$

linear sum
linear correlation.

primary
 $h(x) = [x] \cdot w$...
primary form (data)

• dual form has input sum $\langle x_i, x_j \rangle = K_{ij}$, not the data (x)

$$h(x) = \sum \dots \underbrace{\langle x_j, x_i \rangle}$$

$$\text{sum}(x_i, x_j) = K_{ij}$$

example $x_i = \text{patient in hospital}$

$h(x) = \text{fraction (patient data
blood pressure,
age,
symptoms)}$

dual:
 $h(x) = \text{fraction of similarity
of } (x, \text{other patients})$

Kernel Trick

- $h(x) \Rightarrow$ dual form $h(x) = \text{function}(K_{ij})$
 $\langle x_i, x_j \rangle = \text{sim}(x_i, x_j)$
- $\textcircled{\text{Th}}$ h works with any valid kernel
(not just linear kernel)
- plug in a domain-specific/sim-designed kernel (valid)
ex gaussian kernel $K_{ij} = e^{-\|x_i - x_j\|^2 / \beta}$ (prove valid)
- use the classifier $h(\cdot)$ with non-linear kernel.
- effect: $h(x)$ linear classifier \Rightarrow dual $h(x) = \text{function}$ (non-linear)
can do non-linear separation!
(ex. regression, SVM)

Lecture 7/16 : - KNN classifier - distance/similarities.

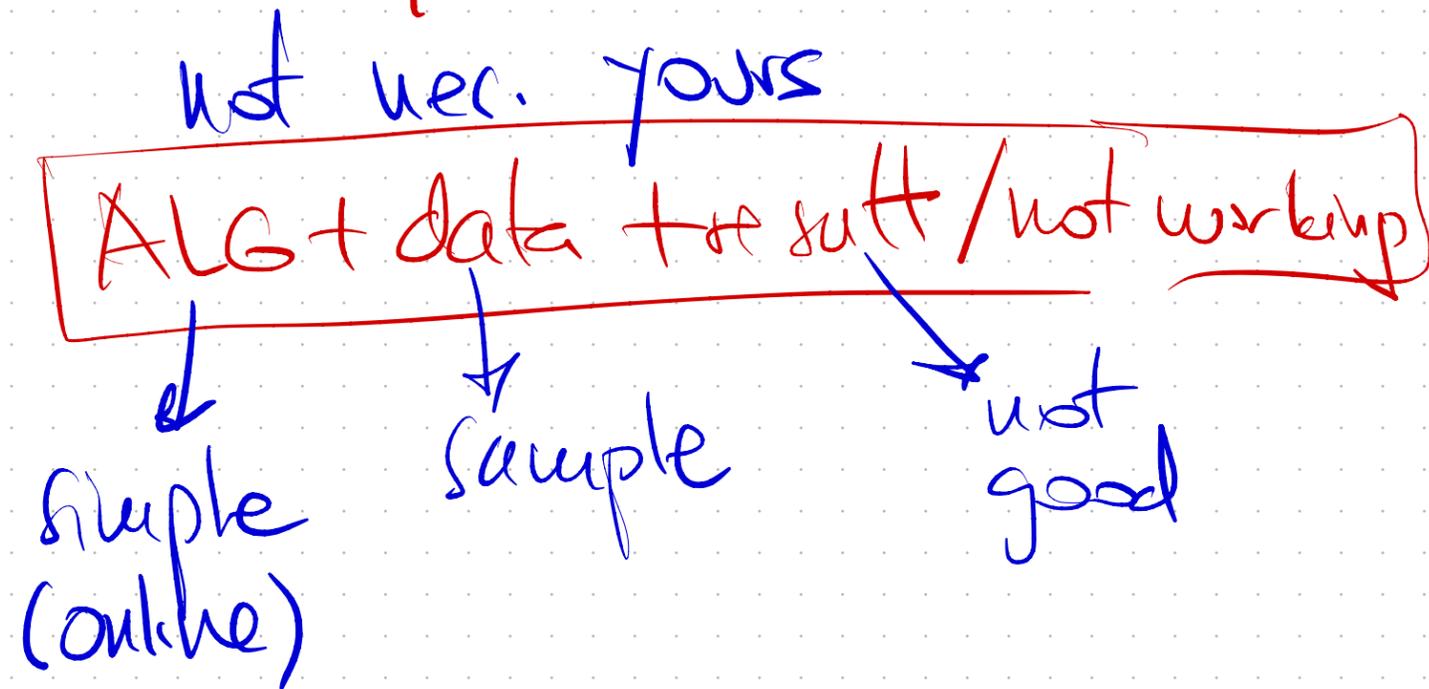
- Kernelization/duality of linear classifiers

= HW 5: SVM + kernels demo Friday 7/18 Caleb

- Projects chat's update

- this week: basic setup

concise
informal



- next week : feedback + direction

• modify ALG

• more data

• diff data

• computation/scale

• debug

• evaluation

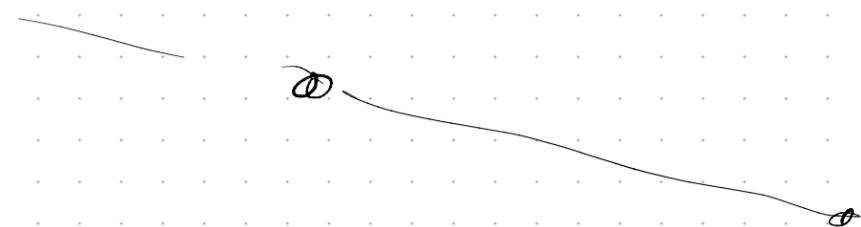
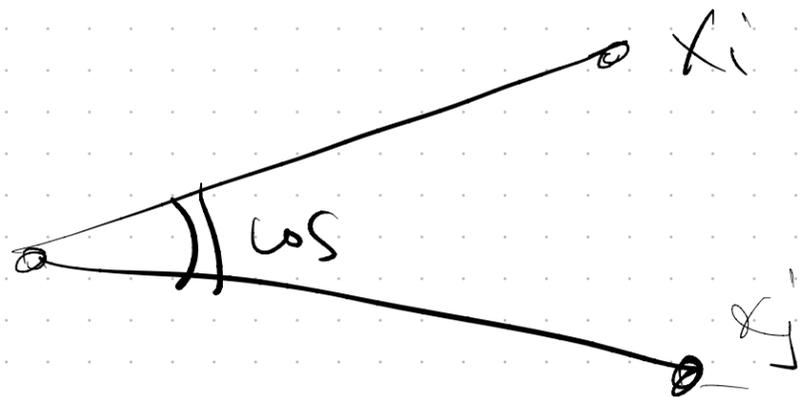
• Other idea.

Fundamental ML: similarity (x_i, x_j) between 2 datapoints
 \approx distance (x_i, x_j) opposite

① correct / exact $\text{sim}(x_i, x_j) \equiv$ correct classification / ML

• default sim with linear algebra = $\langle x_i, x_j \rangle = \sum_{d=1}^D x_i^d \cdot x_j^d$
 \equiv linear kernel $K_{ij} = [X \cdot X^T]_{ij} = \langle x_i, x_j \rangle$

• cosine $(x_i, x_j) = \frac{\langle x_i, x_j \rangle}{\sqrt{\|x_i\| \|x_j\|}} = \frac{\text{dot product}}{\text{normalized by length-given}}$



• Euclidean distance $\|x_i - x_j\|^2 = \sum_{d=1}^D (x_i^d - x_j^d)^2$
 (Norm 2 distance)

• Manhattan distance $\|x_i - x_j\|_1 = \sum_{d=1}^D |x_i^d - x_j^d|$
 Norm 1 distance

• Norm p distance $\|x_i - x_j\|^p = \sum_{d=1}^D |x_i^d - x_j^d|^p$

• Mahalanobis distance $\text{mahal}(x_i, x_j) = \sqrt{(x_i - x_j) \Sigma^{-1} (x_i - x_j)^T}$

• if $\Sigma_{\text{covar}} = \begin{bmatrix} 1 & & 0 \\ & 1 & \\ 0 & & 1 \end{bmatrix}$ \rightarrow unit variance on all feat
 \rightarrow 0 covar between feat \Rightarrow Euclid distance
 $\sqrt{(x_i - x_j) (x_i - x_j)^T}$
 COVAR INVERSE

• If $\Sigma_{\text{covar}} = \begin{bmatrix} \sigma_1^2 & & 0 \\ & \sigma_2^2 & \\ 0 & & \ddots \\ & & & \sigma_D^2 \end{bmatrix}$

\rightarrow no covar between feat
 \Rightarrow weighted / normalized Euclidean Distance
 all features matter same

$$(x_i - x_j) \begin{bmatrix} \frac{1}{\sigma_1^2} & & 0 \\ & \frac{1}{\sigma_2^2} & \\ & & \ddots \\ & & & \frac{1}{\sigma_D^2} \end{bmatrix} (x_i - x_j)^T$$

• Euclidean dist \Rightarrow similarity

Gaussian / RBF kernel $\sigma = \text{scaler}$

$$k_{ij} = \text{sim}(x_i, x_j) = e^{-\frac{\|x_i - x_j\|^2}{2\sigma^2}} = \max k_{ij} = e^0 = 1$$

when $x_i = x_j$

when $k_{ij} = e^{-\infty} = 0$

• edit distance (Hamming)
 $\text{ham}(x_i, x_j) = \sum_{d=1}^D \mathbb{1}[x_i^d \neq x_j^d]$

$\mathbb{1}[x_i^d \neq x_j^d]$ = # of different feat values.
 different $\equiv |x_i^d - x_j^d| \geq \epsilon$

KERNEL = SIMILARITY

$K = [K_{ij}]_{\substack{i=1:N \\ j=1:N}}$ kernel matrix
 $N \times N$

• not any similarity (domain specific) is a valid kernel (math)!

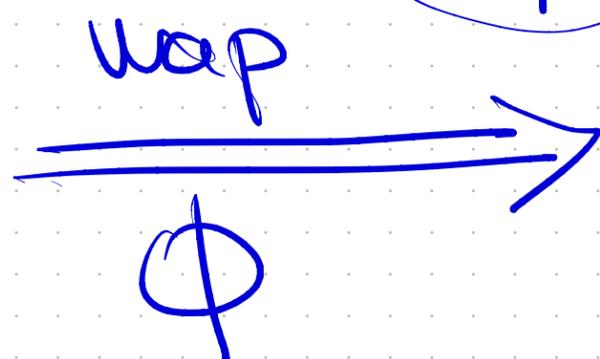
• what are valid kernels? $K =$ always linear but

in implicit space / represent.

X space

ORIGINAL data $D \times d$

$$X = [x^1 \ x^2 \ \dots \ x^D]$$



Φ space

NEW SPACE $G \times d$

$$\Phi(x) = [\phi^1(x) \ \phi^2(x) \ \dots \ \phi^G(x)]$$

$$K = \text{valid kernel} \iff K = \Phi(X) \cdot \Phi(X)^T$$

$$K_{ij} = \langle \Phi(x_i), \Phi(x_j) \rangle$$

linear in Φ -space

• linear kernel $K = X \cdot X^T$ uses no Φ , or $\Phi(x) = x$

Kernel Trick in ML

• don't use Φ explicit, only need existence ($\Phi(x)$) theory
(do not compute $\Phi(x)$ map)
know map

• primal ^{ML} ALG $h(x) = \text{function}(x) \implies$ dual form
 $h(x) = \text{function}(K)$

then we can use $h(x)$ algorithm
with any kernel. (Th, proven math)

k examples

$K_{ij} = (\langle x_i x_j \rangle)^2$ valid kernel. ✓

explicit ϕ for this toy example

original $x = [x^1 \ x^2]$
2-dim

4-dim

$$\phi(x) = [(x^1)^2, (x^2)^2, x^1 \cdot x^2, x^2 \cdot x^1]$$

$$[(x_i^1)^2 \ (x_i^2)^2 \ (x_i^1 x_i^2) \ (x_i^2 x_i^1)] \cdot x$$

$$\langle \phi(x_i) \cdot \phi(x_j) \rangle$$

linear similarity
in ϕ space

$$[(x_j^1)^2 \ (x_j^2)^2 \ (x_j^1 x_j^2) \ (x_j^2 x_j^1)]$$

$$= (x_i^1 \cdot x_j^1)^2 + (x_i^2 x_j^2)^2 + 2 x_i^1 x_j^1 x_i^2 x_j^2$$

$$= (x_i^1 x_j^1 + x_i^2 x_j^2)^2$$

$$= (\langle x_i x_j \rangle)^2$$

calculate w/out ϕ
 $(x_i x_j^T)^2$

$K = (\langle x_i, x_j \rangle + 1)^2$ valid kernel.

calculate $k_{ij} = f(x_i, x_j)$ without ϕ

exercise $k_{ij} = \langle \phi(x_i), \phi(x_j) \rangle$

original

$[x^1, x^2] \rightarrow \phi_{(x)} = [(x^1)^2, (x^2)^2, \sqrt{2}x^1, \sqrt{2}x^2, \sqrt{2}x^1x^2, 1]$ 6-dim

in general: polynomial kernel $k_{ij} = (\langle x_i, x_j \rangle + c)^p$

valid $\Leftrightarrow \exists \phi(x)$ map such that $k_{ij} = \langle \phi(x_i), \phi(x_j) \rangle$

practice: we do not calculate/explicit use ϕ

Gaussian kernel

$$K_{ij} = e^{-\frac{1}{2} \|x_i - x_j\|^2} \quad \underline{\underline{??}} \quad \sum_{t=0}^{\infty} \frac{(x_i x_j^T)^t}{t!} e^{-\frac{1}{2} \|x_i\|^2} e^{-\frac{1}{2} \|x_j\|^2}$$

VALID

$$K = \Phi_{(X)} \Phi_{(X)}^T$$

dot product $\langle \Phi(x_i), \Phi(x_j) \rangle$

Φ infinite dimensionality

calculate w/out Φ : $e^{-\frac{1}{2} \|x_i - x_j\|^2}$

Linear Algorithms

primal \Rightarrow dual $h(x) = \text{function}(x x^T) = \text{function}(K)$

✓ SVM

✓ Linear Reg

• PCA

✓ Perceptron

• KNN with kernels

to enable usage of K kernel

Recall Linear Regression Normal Eq with L_2 -penalty (May 14 Lecture)

intuition $J(w) = \frac{1}{2}(Xw - Y)^T(Xw - Y)$ convex in w

$\Rightarrow \frac{\partial J}{\partial w} = 0$ when $J(w)$ is minimum.

Q?

$$\frac{\partial J}{\partial w} = \frac{\partial}{\partial w} (Xw - Y)^T(Xw - Y)$$

$$= \frac{\partial}{\partial w} (w^T X^T X w - w^T X^T Y - Y^T X w)$$

$$= \frac{\partial}{\partial w} (w^T X^T X w - 2 w^T X^T Y - Y^T Y)$$

$$= \frac{\partial}{\partial w} w^T X^T X w - \frac{\partial}{\partial w} w^T X^T Y$$

$$= \frac{1}{2} 2 X^T X w - X^T Y + \lambda I \cdot w$$

$$= X^T X w - X^T Y + \lambda I w$$

want $= 0$ \rightarrow sq $D \times D$, pos definite (pos eig val)

$$X^T X w = X^T Y \rightarrow (X^T X + \lambda I) w = X^T Y$$

invert-left by $(X^T X)^{-1}$

$$w = (X^T X + \lambda I)^{-1} X^T Y$$

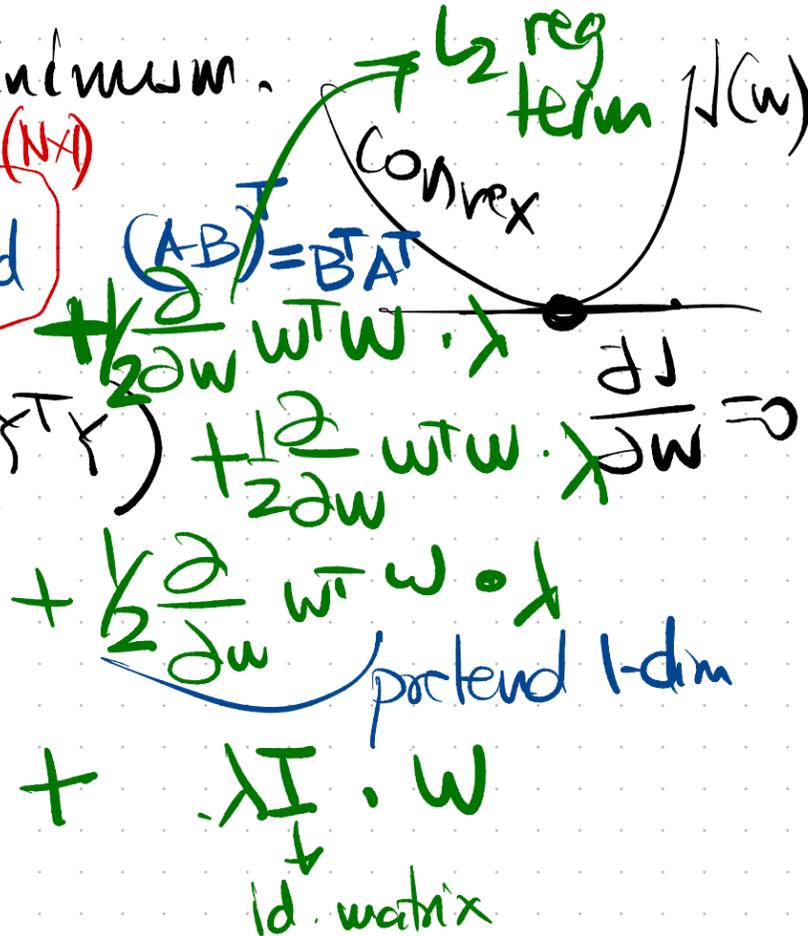
$w = (X^T X + \lambda I)^{-1} X^T Y$ NORMAL EQ (closed form)
Best (in) w for min squares of reg. line

(X) dim $(D) \times (N)$
 (X^T) dim $(N) \times (D)$
same transposed

$w^T X^T Y$

$Y^T X w$

$Y^T Y$
 $\partial = 0$



Dualize Regression.

primal: compute w w/coef

$$[X^T X + \lambda I_D] w = X^T y$$

id $D \times D$

dual: $w = X^T \alpha$
var α = lin comb (X^T)
with coef α .

$$\alpha = [X X^T + \lambda I_N]^{-1} Y$$

$N \times 1$ $N \times N$ $N \times N$ $N \times 1$

α_i = coef for x_i

proof that $w = X^T \alpha$

check the calculation assuming $w = X^T \alpha$

eg $(X^T X + \lambda I_D) w \stackrel{?}{=} X^T y$

$$[X^T X + \lambda I_D] \cdot X^T \alpha = X^T X X^T \alpha + \lambda X^T \alpha$$

$$= X^T (X X^T \alpha + \lambda I_N \alpha) = X^T (X X^T + \lambda I_N) \alpha$$

$$= X^T [X X^T + \lambda I_N] [X X^T + \lambda I_N]^{-1} Y$$

same

$$= X^T y \quad \checkmark$$

$$[X X^T + \lambda I_N]^{-1} Y$$

dual L_2 -regression (Kernelized Ridge Regression, HW 5)

$$\alpha_i = \text{dual var} \quad w = x^T \alpha$$

$i=1:N$

$$h(x) = x \cdot w = x \cdot \boxed{x^T \alpha} = \boxed{x \cdot x^T} \alpha = K \alpha$$

primal

$$h(z) = \sum_{i=1}^N \alpha_i K(x_i, z)$$

= function (similarities)

test for datapoint z

train

$$\alpha = [xx^T + \lambda I_n]^{-1} y$$

use K = kernel valid other than xx^T

train

$$\alpha = [K + \lambda I_n]^{-1} y$$

test

$$h(z) = \sum \alpha_i K(z, x_i)$$

K = Gaussian, or K = polynomial

Recall Perceptron update rule (May 28 Lecture)

Perceptron problem (modified data)

- Given N vectors x_1, x_2, \dots, x_N

- Find hyperplane w s.t. $x_i w \geq 0 \forall i$
(puts all points on pos side)

error $J(w) = \sum_{x_i} -x_i w \geq 0$

mistake
 $x_i w < 0$

pos
how far from
hyperplane x_i is

not mistake
(correct) $x_i w \geq 0$

OBS: no error \Rightarrow no mistakes
 $\Rightarrow J(w) = 0$

$x_i w < 0$

update: Every mistake x_i gets added to $w =$ normal of hyperplane

Gradient Desc update

$$\frac{\partial J}{\partial w^d} = - \frac{\partial x_i w}{\partial w^d} = - \frac{\partial [\sum_d x_i^d w^d]}{\partial w^d}$$

$$= -x_i^d \quad \text{if } x_i =$$

$$\frac{\partial J}{\partial w} = -x_i^T = \text{mistake}$$

Vector
GD update for datapoint x_i

$$w_{\text{new}} = w_{\text{old}} + \lambda \cdot x_i^T$$

(all datapoints $x_i =$ mistakes)

$$w_{\text{new}} = w_{\text{old}} + \lambda \sum_{x_i \text{ mistakes}} x_i^T$$

$\{x_i w < 0\}$

Perceptron Algorithm

while (condition)
step k : x_k misclassified
 $w^{\text{new}} = w^{\text{old}} + x_k$

mechanics: count $\alpha_i = \#$ steps where x_i has been used for update.

$$w_{\text{final}} = 0 + \sum_{i=1}^N \alpha_i \cdot x_i = x^T \cdot \alpha$$

add to w

counts for many times each datapoint

$$= \alpha_1 x_1 + \alpha_2 x_2 + \dots + \alpha_N x_N$$

#times x_1 has been added to w

$$h(x) = x \cdot w = x \cdot (x^T \alpha) = K \cdot \alpha$$

use any valid kernel!