



gets very slow very fast



## The kernel Trick

recall polynomial features: w/  $d$  input features degree- $p$  polynomials blow up to  $O(d^p)$  features

e.g.  $p=4$ , features blown up to  $d^4$

Today, we use those features w/o computing them

key observations: In many algorithms,

- weights are linear combo of the sample points

- can use inner products of lifted vectors  $\vec{z}(x)$  only - don't need each  $\vec{z}(x)$  explicitly computed

$$\vec{w} = \vec{x}^T \vec{a} = \sum_{i=1}^n a_i x_i, \vec{a} \in \mathbb{R}^n$$

each comp. of  $\vec{a}$  is  
a dual weight

If this is true for an algorithm, plug in  $\vec{w} = \vec{x}^T \vec{a}$  into the algo.

↳ get rid of  $\vec{w}$ , optimize the  $n$  dual weights  $\vec{a}$  (dual parameters) instead of the  $d+1$  (or  $d^p$  if polynomial features) primal weights of  $\vec{w}$

## Kernel Ridge Regression

problem: impossible to write weights as lin. combo. of data if we have a bias term

↳ 1 needs to penalize bias too

↳ to mitigate this effect, center  $X$  &  $\vec{y}$  → mean of 0

$$x'_i \leftarrow x_i - \mu_x, y'_i \leftarrow y_i - \mu_y, x'_{i,d+1} = 1$$

col. of 1's appended

- This lets us replace  $I'$  w/  $I$  in normal equations:

$$(X^T X + 1I) \vec{w} = X^T \vec{y}$$

Primal

↳  $I'$  is w/ bias term

Suppose  $\vec{a}$  is a solution to:  $(X^T X + 1I) \vec{a} = \vec{y}$  dual

$$\Rightarrow X^T \vec{y} = X^T X \vec{a} + 1 \vec{a} = (X^T X + 1I) \vec{a}$$

lin. comb. of  
sample pts.

(dual solution)  $\Rightarrow \vec{a}$  sol as long as  $\vec{y} = X^T \vec{a}$

$$\text{Primal: } \|X\vec{v} - \vec{y}\|^2 + \lambda \|\vec{v}\|^2$$



$$\text{dual form of ridge regression: } \min_{\vec{a}} \|X^T \vec{a} - \vec{y}\|^2 + \lambda \|\vec{a}\|^2$$

training part of algo: Solve  $(X^T X + \lambda I)\vec{a} = \vec{y}$  for  $\vec{a}$  ✓ big if using listed features

$$\text{Testing: Regression function is: } h(\vec{z}) = \vec{w}^T \vec{z} = \vec{a}^T \underbrace{[X \vec{z}]}_n = \sum_{i=1}^n a_i (\vec{x}_i^T \vec{z})$$

Weighted sum of inner products

define  $k(\vec{x}, \vec{z})$  as the kernel fn.

$$\text{For now, } k(\vec{x}, \vec{z}) = \vec{x}^T \vec{z}. \text{ later, } \vec{x} = \Phi(\vec{x}), \vec{z} = \Phi(\vec{z})$$

let  $K = X X^T$  be the  $n \times n$  kernel matrix,  $k_{ij} = k(x_i, x_j)$  always PSD

Singular  $K$  if  $n > d+1$  (Sometimes still Singular even if this isn't the case)  
↳ full  $X$

↳ almost always true w/ large datasets  
↳ choose  $\lambda > 0$  to solve system

↑ if no kernel trick

dual ridge regression algorithm:

training {  
1 Compute the kernel matrix one entry at a time:  $K_{ij} \leftarrow k(x_i, x_j)$   $\Theta(n^2 d)$  time

2 Solve  $(K + \lambda I)\vec{a} = \vec{y}$  for  $\vec{a}$   $\Theta(n^3)$  time

testing {  
3 for each test point  $\vec{z}$ :

$$h(\vec{z}) = \sum_{i=1}^n a_i k(x_i, \vec{z})$$

\* again, same sol. as primal ridge reg.

algo (if bias penalized)

Notice: Sample pts. only used as inner products

Dual: Solve  $n \times n$  linear system  $\Theta(n^3 + n^2 d)$  time

preferred when  $d > n$

Primal: Solve  $d \times d$  linear system  $\Theta(d^3 + d^2 n)$  time

advantage of dual: When using polynomial features,  $d$  in primal algo increases but  $d$  in dual won't

The kernel Trick / kernelization: We can compute a polynomial kernel that involves many monomial terms w/o ever computing those terms

Polynomial kernel of degree p:  $k(\vec{x}, \vec{z}) = (\vec{x}^T \vec{z} + 1)^p$

Theorem:  $(\vec{x}^T \vec{z} + 1)^p = \Phi(\vec{x})^T \Phi(\vec{z})$ ,  $\Phi(\vec{x})$  is lifting fn that contains every monomial term w.r.t.  $\vec{x}$  of degree 0..p

e.g. d=2, p=2

$$\rightarrow (\vec{x}^T \vec{z} + 1) = x_1^2 z_1^2 + x_2^2 z_2^2 + 2x_1 z_1 x_2 z_2 + 2x_1^2 + 2x_2 z_2 + 1$$

$\Phi(\vec{x})$   $\Phi(\vec{z})$

$$= \begin{bmatrix} x_1^2 & z_1^2 & \sqrt{2} x_1 x_2 & \sqrt{2} x_1 & \sqrt{2} x_2 & 1 \end{bmatrix} \begin{bmatrix} z_1^2 \\ z_2^2 \\ \sqrt{2} z_1 z_2 \\ \sqrt{2} z_1 \\ \sqrt{2} z_2 \\ 1 \end{bmatrix}$$

(can't control the constants, but that's okay)

$$= \Phi(\vec{x})^T \Phi(\vec{z})$$

When optimized, primal weights will automatically adjust themselves to compensate for the constants

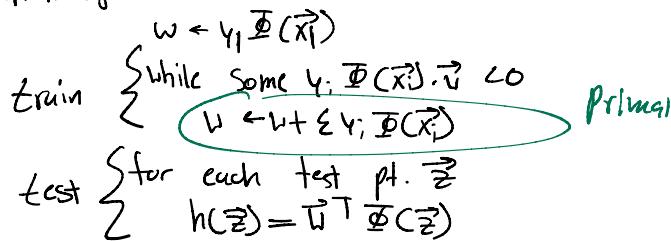
key win: compute  $\Phi(\vec{x})^T \Phi(\vec{z})$  in  $O(d)$  time instead of  $O(p)$  time even though  $\Phi(\vec{x})$  &  $\Phi(\vec{z})$  have length  $O(d^p)$

kernel ridge reg. w/ polynomial features: replace  $x_i$  w/  $\Phi(x_i)$ :

let  $k(\vec{x}, \vec{z}) = \underbrace{\Phi(\vec{x})^T \Phi(\vec{z})}_{\text{NEVER COMPUTED}} = (\vec{x}^T \vec{z} + 1)^p$

## kernel Perceptrons:

Facilitated perceptron algo:



Let  $\Phi(X)$  be the  $n \times D$  matrix w/ rows  $\Phi(\vec{x}_i)^T$ ,  $D = \text{length of } \Phi(\vec{x}_i)$

$$\Rightarrow k = \Phi(X) \Phi(X)^T$$

→ Dualize w/  $\vec{w} = \Phi(X)^T \vec{a}$

↳ Then the code " $a_i \leftarrow a_i + \varepsilon y_i$ "  $\equiv \vec{w} \leftarrow \vec{w} + \varepsilon y_i \Phi(\vec{x}_i)$   
also,  $\vec{w} = \mathcal{I}(X)^T \vec{a}$

$$\Rightarrow \mathcal{I}(x_i)^T \vec{w} = (\mathcal{I}(X) \vec{a})_i - (\mathcal{I}(X) \mathcal{I}(X)^T \vec{a})_i = (k \vec{a})_i$$

⇒ Dual Perceptron Algo:

- 1  $\vec{a} \leftarrow [y_1 \ 0 \ \dots \ 0]^T$
- 2 Compute the kernel matrix  $K_{ij} \leftarrow k(\vec{x}_i, \vec{x}_j)$   $\mathcal{O}(n^2 d)$  time (kernel init)
- 3 **while** some  $y_i(K\vec{a})_i < 0$ :
- 4      $a_i \leftarrow a_i + \varepsilon y_i$   $\mathcal{O}(1)$
- 5
- 6 **for** each test point  $\vec{z}$
- 7      $h(\vec{z}) \leftarrow \sum_{j=1}^n a_j k(\vec{x}_j, \vec{z})$   $\mathcal{O}(nd)$  time (kernel Z init)

## Kernel Logistic Regression

Stochastic gradient descent Step:

$$a_i \leftarrow a_i + \epsilon (y_i - S((ka)_i))$$

Sigmoid

Batch gradient descent:  $a \leftarrow a + \epsilon (Y - S(ka))$

applied component-wise to vector

for each test pt.  $\vec{z}$

$$h(\vec{z}) = S\left(\sum_{j=1}^n a_j k(x_j, \vec{z})\right) \quad S(\cdot) \text{ for posterior probabilities}$$

Gaussian kernel: Use feature vectors in an infinite-dimensional Space

Gaussian kernel based on observation that  $\exists \Phi(\vec{x})$  such that

$$k(\vec{x}, \vec{z}) = \exp\left(-\frac{\|\vec{x} - \vec{z}\|^2}{2\sigma^2}\right) = \Phi(\vec{x})^\top \Phi(\vec{z})$$

isotropic, centred at  $\vec{z}/\vec{x}$

e.g.  $d=1$

$$\Phi(\vec{x}) = \exp\left(-\frac{x^2}{2\sigma^2}\right) \left[ 1 \ \frac{x}{\sigma\sqrt{1}} \ \frac{x^2}{\sigma^2\sqrt{2!}} \ \dots \right]^\top$$

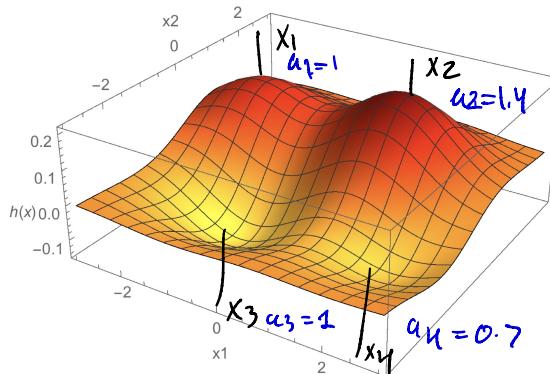
$$\Phi(\vec{x})^\top \Phi(\vec{z}) = k(\vec{x}, \vec{z})$$

more helpful to think in terms of fn  $k(\cdot)$  rather than  $\Phi(\cdot)$

Key observation: hypothesis  $h(\vec{z}) = \sum_{j=1}^n a_j k(x_j, \vec{z})$  is a lin. combo of Gaussians  
dual weights are coefficients  
Centred at sample pts.  
basis for hypothesis

4 Sample pts.

- height dictated by coeff.  $\alpha_i$



Gaussian kernel popular in practice, but why?

- gives very smooth regression fn.  $h(\cdot)$ ; infinitely differentiable,  $C^\infty$  continuous
- $h(\cdot)$  behaves somewhat like k-nearest-neighbors, but smoother
- oscillates less than polynomials (depending on  $\sigma$ )
- $k(\vec{x}, \vec{z})$  can be interpreted as a similarity measure
  - ↳ max when  $\vec{z} = \vec{x}$
  - ↳ goes to 0 as distance increases
- Sample pts. "vote" for value at  $\vec{z}$ , but closer pts. get more weight of vote

$\sigma$  is an h-param chosen through validation, trades off bias vs. variance

- larger  $\sigma \rightarrow$  wider Gaussians  $\rightarrow$  smoother  $h(\cdot)$   $\rightarrow$  more bias (fn. can't adapt to data)  
 $\rightarrow$  less variance