

HW 2 is due

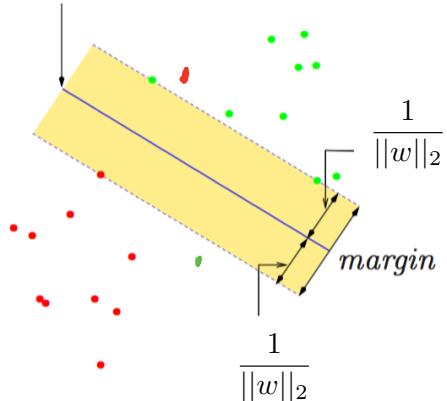
Kernels



W

What if the data is not linearly separable?

$$\bullet x^T w + b = 0$$



some points don't satisfy margin constraint:

$$\min_{w,b} ||w||_2^2$$

$$y_i(x_i^T w + b) \geq 1 \quad \forall i$$

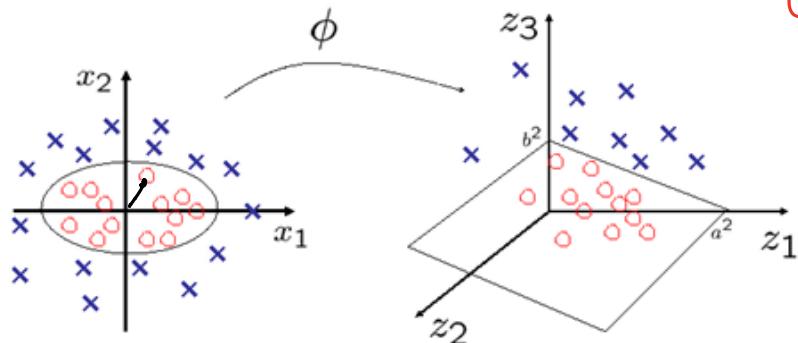
Two options:

1. Introduce slack to this optimization problem
2. Lift to higher dimensional space

$$x \mapsto \phi(x)$$

What if the data is not linearly separable?

Use features of features of features...



$$(x_1, x_2) \rightarrow (r, \theta)$$
$$r = \sqrt{x_1^2 + x_2^2}$$
$$\theta = (\delta) \frac{x_1 \cdot x_2}{\|x_1\| \|x_2\|}$$

$$\begin{pmatrix} z_1 \\ z_2 \\ z_3 \end{pmatrix} = \begin{pmatrix} x_1 \\ x_2 \\ \sqrt{x_1^2 + x_2^2} \end{pmatrix}$$

Creating Features

$$x \in \mathbb{R}^d, \quad h(x) \in \mathbb{R}^p$$

Transformed data:

$h : \mathbb{R}^d \rightarrow \mathbb{R}^p$ maps original features to a rich, possibly high-dimensional space

for $d > 1$, generate

$$\{u_j\}_{j=1}^p \subset \mathbb{R}^d$$

$$h_j(x) = (u_j^T x)^2$$

$$h_j(x) = \frac{1}{1 + \exp(u_j^T x)}$$

$$h_j(x) = \cos(u_j^T x)$$

u_j ; random

in $d=1$:

$$h(x) = \begin{bmatrix} h_1(x) \\ h_2(x) \\ \vdots \\ h_p(x) \end{bmatrix} = \begin{bmatrix} x \\ x^2 \\ \vdots \\ x^p \end{bmatrix}$$

Creating Features

Transformed data:

ϕ : $h : \mathbb{R}^d \rightarrow \mathbb{R}^p$ maps original features to a rich, possibly high-dimensional space

for d>1, generate

$$\text{in } d=1: \quad h(x) = \begin{bmatrix} h_1(x) \\ h_2(x) \\ \vdots \\ h_p(x) \end{bmatrix} = \begin{bmatrix} x \\ x^2 \\ \vdots \\ x^p \end{bmatrix}$$
$$\{u_j\}_{j=1}^p \subset \mathbb{R}^d$$
$$h_j(x) = (u_j^T x)^2$$
$$h_j(x) = \frac{1}{1 + \exp(u_j^T x)}$$
$$h_j(x) = \cos(u_j^T x)$$

Feature space can get really large really quickly!

Degree-d Polynomials

input: $u \in \mathbb{P}^2$, (u_1, u_2)

degree-1 : $\phi(u) = \begin{pmatrix} u_1 \\ u_2 \end{pmatrix} \in \mathbb{P}^2$

degree-2 : $\phi(u) = \begin{pmatrix} u_1 \cdot u_1 \\ u_1 \cdot u_2 \\ u_2 \cdot u_1 \\ u_2 \cdot u_2 \end{pmatrix} \in \mathbb{P}^4$

degree-3 :

$$\phi(u) = \begin{pmatrix} u_1 \cdot u_1 \cdot u_1 \\ u_1 \cdot u_1 \cdot u_2 \\ u_1 \cdot u_2 \cdot u_1 \\ u_1 \cdot u_2 \cdot u_2 \\ u_2 \cdot u_1 \cdot u_1 \\ u_2 \cdot u_1 \cdot u_2 \\ u_2 \cdot u_2 \cdot u_1 \\ u_2 \cdot u_2 \cdot u_2 \end{pmatrix} \in \mathbb{P}^8$$

if $u \in \mathbb{P}^p$
 $\phi(u) = 0 \quad (p^d)^{\dim}$

exponential p^d
 $\Rightarrow n: \# \text{ of data}$

How do we deal with high-dimensional lifts/data?

A fundamental trick in ML: use kernels

A function $\underline{K : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}}$ is a kernel for a map ϕ if $\underline{K(x, x')} = \phi(\underline{x}) \cdot \phi(\underline{x'})$ for all x, x' .

So, if we can represent our algorithms/decision rules as dot products and we can find a kernel for our feature map then we can avoid explicitly dealing with $\phi(x)$.

Linear Regression as Kernels

Training

solution
regularized
least square

$$X \in \mathbb{R}^{n \times d}, Y \in \mathbb{R}^{n \times 1}$$

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

$$= \frac{1}{d \times n} X^T (X X^T + \lambda I_n)^{-1} Y$$

$$X = \begin{bmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_n^T \end{bmatrix}$$

(linear algebra operations, e.g. SVD)

Decision rule
prediction

$$x_{\text{new}} \in \mathbb{R}^d$$

$$\hat{y}_{\text{new}} = \hat{w}^T x_{\text{new}}$$

$$= y^T (X X^T + \lambda I_n)^{-1} \underbrace{x_{\text{new}}}_{|X|}$$

$$x_{\text{new}} = \begin{bmatrix} x_1^T \\ \vdots \\ x_n^T \end{bmatrix} \begin{bmatrix} x_{\text{new}} \end{bmatrix}_d$$

$$X X^T \in \mathbb{R}^{n \times n}$$

$$\stackrel{i,j}{\in \{1, \dots, n\}} [X X^T]_{i,j} = \langle x_i, x_j \rangle$$

$$[X x_{\text{new}}]_i = \langle x_i, x_{\text{new}} \rangle$$

$$\Rightarrow \langle \cdot, \cdot \rangle \rightarrow k(\cdot, \cdot)$$

$$\Rightarrow k(x_i, x_j)$$

$$k(x_i, x_{\text{new}})$$

Dot-product of polynomials

$$u \in \mathbb{P}^2, \quad v \in \mathbb{P}^2$$

$$\begin{aligned} u &\in \mathbb{P}^3, \quad d=2 \\ (u_1, u_2, u_3)^\top & (u_1, u_2, u_3) \\ \Rightarrow q &= 3^2 \end{aligned}$$

$\Phi(u) \cdot \Phi(v) = \text{polynomials of degree exactly } d$

$$d=1 : \phi(u) = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} \quad \langle \phi(u), \phi(v) \rangle = u_1 v_1 + u_2 v_2$$

$$d=2 : \phi(u) = \begin{bmatrix} u_1^2 \\ u_2^2 \\ u_1 u_2 \\ u_2 u_1 \end{bmatrix} \quad \langle \phi(u), \phi(v) \rangle = \underbrace{u_1^2 v_1^2 + u_2^2 v_2^2}_{K(u,v)} + 2u_1 u_2 v_1 v_2$$

degree - d polynomial with $u \in \mathbb{P}^p$
 if use explicit feature, # of operations $O(p^d)$

$K(u, v) = (\langle u, v \rangle)^d$, # of operations,
 $O(p + \log d)$

Dot-product of polynomials

$\Phi(\mathbf{u}) \cdot \Phi(\mathbf{v}) = \text{polynomials of degree exactly } d$

$$d=1 : \phi(u) = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} \quad \langle \phi(u), \phi(v) \rangle = u_1 v_1 + u_2 v_2$$

$$d=2 : \phi(u) = \begin{bmatrix} u_1^2 \\ u_2^2 \\ u_1 u_2 \\ u_2 u_1 \end{bmatrix} \quad \langle \phi(u), \phi(v) \rangle = u_1^2 v_1^2 + u_2^2 v_2^2 + 2u_1 u_2 v_1 v_2$$

Feature space can get really large really quickly!

General d : Dimension of $\phi(u)$ is roughly p^d if $u \in \mathbb{R}^p$

Feature expansion can be written **implicitly** $K(\mathbf{u}, \mathbf{v}) = (\mathbf{u} \cdot \mathbf{v})^{\frac{1}{p}}$

δ, ν : hyper-parameter

Examples of Kernels

- Polynomials of degree exactly d

$$K(\mathbf{u}, \mathbf{v}) = (\mathbf{u} \cdot \mathbf{v})^{\underline{d}}$$

- Polynomials of degree up to d

$$K(\mathbf{u}, \mathbf{v}) = (\mathbf{u} \cdot \mathbf{v} + 1)^{\underline{d}}$$

- Gaussian (squared exponential) kernel

$$K(\mathbf{u}, \mathbf{v}) = \exp\left(-\frac{\|\mathbf{u} - \mathbf{v}\|^2}{2\sigma^2}\right)$$

- Sigmoid

$$K(u, v) = \tanh(\gamma \cdot u^T v + r)$$

(RBF Kernel)

tanh(x)
 $= \frac{e^x - e^{-x}}{e^x + e^{-x}}$

The Kernel Trick

Data $\{(x_i, y_i)\}_{i=1}^n$

(1) Pick a kernel $K : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$

α can depend
on $\langle x_i, x_j \rangle$
and y

(2) For a linear predictor, show $w = \sum_{i=1}^n \alpha_i x_i$, $\alpha \in \mathbb{R}$

(3) Change loss function/decision rule to only access data through dot products

$$w^\top x_{\text{new}} = \sum_{i=1}^n \alpha_i x_i^\top x_{\text{new}}$$

Substitute $K(x_i, x_j)$ for $x_i^\top x_j$

$$\hat{y}_{\text{new}} = \sum_{i=1}^n \alpha_i K(x_i, x_{\text{new}})$$

Loss Functions

$$\{(x_i, y_i)\}_{i=1}^n \quad x_i \in \mathbb{R}^d \quad y_i \in \mathbb{R}$$

- Loss functions:

$$\sum_{i=1}^n \ell_i(w)$$

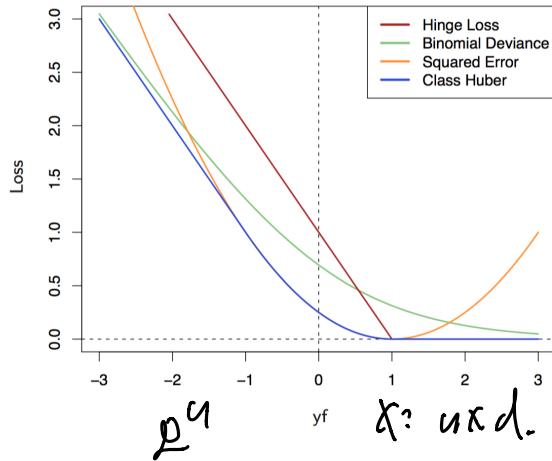
Squared error Loss: $\ell_i(w) = (y_i - x_i^T w)^2$

Logistic Loss: $\ell_i(w) = \log(1 + \exp(-y_i x_i^T w))$

0/1 loss: $\ell_i(w) = \mathbb{I}[\text{sign}(y_i) \neq \text{sign}(x_i^T w)]$

Hinge Loss: $\ell_i(w) = \max\{0, 1 - y_i x_i^T w\}$

SVM



$$\begin{aligned}w &= \alpha^T X \\x_i^T w &= \alpha^T X X_i \\&= \sum_{j=1}^n \alpha_j \langle x_i, x_j \rangle\end{aligned}$$

The Kernel Trick for regularized least squares

$$\alpha \in \mathbb{R}^n$$

$$\hat{w} = \arg \min_w \sum_{i=1}^n (y_i - x_i^T w)^2 + \lambda \|w\|_2^2$$

There exists an $\alpha \in \mathbb{R}^n$: $\hat{w} = \sum_{i=1}^n \alpha_i x_i$

$$\hat{\alpha} = \underset{\alpha}{\operatorname{argmin}} \sum_{i=1}^n \left(y_i - \sum_{j=1}^n \alpha_j x_i^T x_j \right)^2 + \lambda \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j x_i^T x_j$$

$$(\langle x_i, x_j \rangle \rightarrow K(x_i, x_j))$$

$$\Rightarrow \underset{\alpha}{\operatorname{argmin}} \sum_{i=1}^n \left(y_i - \sum_{j=1}^n \alpha_j K(x_i, x_j) \right)^2 + \lambda \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j K(x_i, x_j)$$

(Quadratic form)

$$\underset{\alpha}{\operatorname{argmin}} \|y - K\alpha\|_2^2 + \lambda \alpha^T K \alpha$$

$K \in \mathbb{R}^{n \times n}, [K]_{ij} = K(x_i, x_j)$

$$\begin{aligned} [\bar{y} - K\bar{\alpha}]_i &= y_i - \bar{K}\bar{\alpha}_i \\ [\bar{K}\bar{\alpha}]_i &= \sum_{j=1}^n \bar{\alpha}_j K(x_i, x_j) \\ [\bar{K}]_{ij} &= \frac{1}{n} \sum_{k=1}^n K(x_i, x_k) - \frac{1}{n^2} \sum_{k=1}^n \sum_{l=1}^n K(x_k, x_l) \end{aligned}$$

The Kernel Trick for regularized least squares

$$\hat{w} = \arg \min_w \sum_{i=1}^n (y_i - x_i^T w)^2 + \lambda ||w||_w^2$$

There exists an $\alpha \in \mathbb{R}^n$: $\hat{w} = \sum_{i=1}^n \alpha_i x_i$

$$\begin{aligned}\hat{\alpha} &= \arg \min_{\alpha} \sum_{i=1}^n (y_i - \sum_{j=1}^n \alpha_j \langle x_j, x_i \rangle)^2 + \lambda \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j \langle x_i, x_j \rangle \\ &= \arg \min_{\alpha} \sum_{i=1}^n (y_i - \sum_{j=1}^n \alpha_j K(x_i, x_j))^2 + \lambda \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j K(x_i, x_j) \\ &= \arg \min_{\alpha} \|\mathbf{y} - \mathbf{K}\alpha\|_2^2 + \lambda \alpha^T \mathbf{K} \alpha\end{aligned}$$

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