# 9/19/2023: Kernels Continued

$K(x, x')$ measures similarity between $x$ & $x'$

↳ inputs to model

similar pairs $(x, x')$ should have large $K(x, x')$

Idea: Make prediction by computing $\sum_{i=1}^{n} \alpha_i K(x^{(i)}, x^{test})$

If $> 0$ → predict 1
If $< 0$ predict $-1$

learned parameter

$i$th training example

test example

Logistic regression already does this (in a way) if you define

$$K(x, x') = x^T x'$$

## Original L.R

Training: $w^{(0)} \leftarrow \underline{0}$

$$w^{(t)} \leftarrow w^{(t-1)} + \eta \cdot \frac{1}{n} \sum_{i=1}^{n} \underbrace{\sigma(-y^{(i)} w^{(t-1)T} x^{(i)}) \cdot y^{(i)}}_{\text{Scalar}} \cdot x^{(i)}$$

Test time: given $x^{test}$

Compute $w^{(final)T} x^{test}$

## Kernelized L.R

Define $w = \sum_{c=1}^{n} \alpha_i x^{(i)}$

Training: $\alpha^{(0)} \leftarrow \underline{0}$

$$\alpha_i^{(t)} \leftarrow \alpha_i^{(t-1)} + \eta \cdot \frac{1}{n} \cdot \sigma(-y^{(i)} \boxed{w^{(t-1)T} x^{(i)}}) \cdot y^{(i)}$$

Do this for $i = 1, ..., n$

$$= \sum_{j=1}^{n} \alpha_j^{(t-1)} \boxed{x^{(j)T} x^{(i)}}$$
$$= K(x^{(j)}, x^{(i)})$$

Test time:

Compute $\sum_{i=1}^{n} \alpha_i^{(final)} x^{(i)T} x^{test}$

## Why do this?

Run kernelized algorithm replacing dot products with other $K(x, x')$

hyperparameter you choose

# Kernels & Features

| y | $x_1$ | $x_2$ |
|---|---|---|
| +1 | 2 | 3 |
| -1 | 0 | 1 |
| ⋮ | | |

transform each row →

| y | $\sqrt{2}\,x_1$ | $\sqrt{2}\,x_2$ | 1 | $x_1^2$ | $x_2^2$ | $\sqrt{2}\,x_1 x_2$ |
|---|---|---|---|---|---|---|
| +1 | $2\sqrt{2}$ | $3\sqrt{2}$ | 1 | 4 | 9 | $6\sqrt{2}$ |
| -1 | 0 | $1\sqrt{2}$ | 1 | 0 | 1 | 0 |
| ⋮ | | | | | | |

Think of this as function $\Phi : \mathbb{R}^2 \to \mathbb{R}^6$

Drawback: using $\Phi$ + logistic regression is $\approx 3x$ slower

For some $\Phi$, you can quickly compute

$$K(x, x') = \Phi(x)^T \Phi(x')^T$$

without actually computing $\Phi(x)$ or $\Phi(x')$

"Kernel trick"

E.g. Quadratic Kernel:

$$K(x, x') = \underbrace{(x^T x' + 1)^2}_{} = \Phi(x)^T \Phi(x')$$

only requires Operations in original vector space $\mathbb{R}^d$

where $\Phi(x) = \begin{bmatrix} 1 \\ \sqrt{2}\,x_1 \\ \vdots \\ \sqrt{2}\,x_d \\ x_1^2 \\ \vdots \\ x_d^2 \\ \sqrt{2}\,x_1 x_2 \\ \vdots \\ \sqrt{2}\,x_{d-1} x_d \end{bmatrix}$

} constant term

} All linear terms $\times \sqrt{2}$

} All $x_i^2$ terms

} All $x_i x_j$ terms $\times \sqrt{2}$

vector is size $O(d^2)$

More Generally: for degree $p$,

$$K(x, x') = (x^T x' + 1)^p = \phi(x)^T \phi(x')$$

for some $\phi$ that includes all monomials of degree $\leq p$

What about RBF?

Fact:

$$\exp\left(\frac{\|x - x'\|^2}{2\sigma^2}\right) = \phi(x)^T \phi(x')$$

for some $\phi(x)$ that is infinite-dimensional

Runtime comparisons: Let's use polynomial kernel of degree $p$

| Original C.R. | Kernelized C.R. |
|---|---|

Original C.R.
- Map each $x^{(i)}$ to $O(d^p)$-size feature vector
- Training: 1 iteration takes $O(n d^p)$
- Testing: takes $O(d^p)$

Better dependence on $n$

Kernelized C.R.
- Use kernel trick
- Training: 1 iteration takes $O(n^2 d)$
- Testing: compute
$$\sum_{i=1}^{n} \alpha_i K(x^{(i)}, x^{test})$$
takes $O(nd)$

Better dependence on $d$ & $p$

Bad if dataset is very large

# Support vector machines (SVM)

Similarities to Logistic Regression
- Binary classification
- Learn linear decision boundary
- Parameter is $w \in \mathbb{R}^d$
- Decision boundary defined by $w^T x = 0$
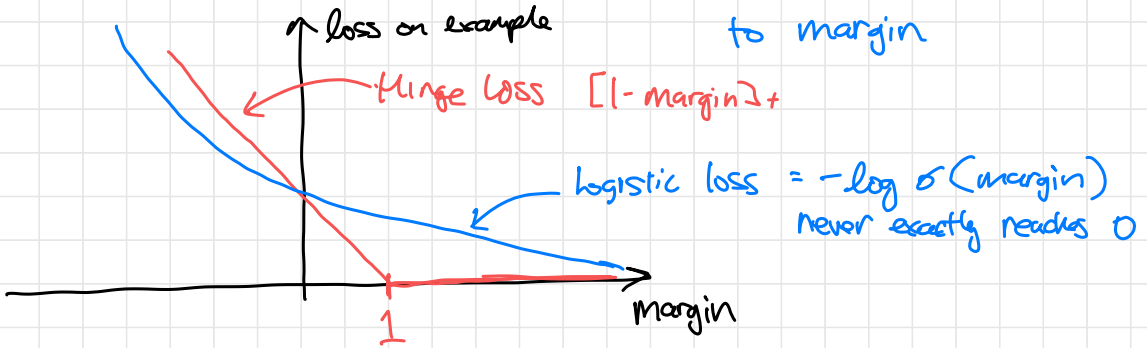
Difference: SVM has no probabilistic interpretation

SVM minimizes the following loss:

$$L(w) = \frac{1}{n} \sum_{i=1}^{n} \left[ 1 - y^{(i)} w^T x^{(i)} \right]_+ \quad + \quad \underbrace{\lambda \|w\|^2}_{L_2 \text{ regularization}}$$

where the margin is $\underbrace{1 - y^{(i)} w^T x^{(i)}}_{\text{margin}}$

Applies "hinge loss" function $f(z) = [1 - z]_+$ to margin

where $[z]_+ = \begin{cases} z & \text{if } z > 0 \\ 0 & \text{if } z \le 0 \end{cases}$

AKA $[z]_+ = \max(z, 0)$

loss on example

Hinge loss $[1 - \text{margin}]_+$

logistic loss $= -\log \sigma(\text{margin})$
never exactly reaches 0

margin

1

What does hinge loss do?

+ unimportant

not-ideal boundary

$w^T x = 1$

SVM learns this
$w^T x = 0$

$w^T x = -1$

distance $\frac{1}{\|w\|}$

SVM wants:
- 0 hinge loss
- small $\|w\|$
$\iff$
large $\frac{1}{\|w\|}$
$\iff$
large distance between examples & decision boundary

Certain examples are _Support vectors_ (green)
in particular, ones where margin $\leq 1$

Ideal $w$ only depends on support vectors
i.e. it is a linear combination of Support vectors only

## Connection to Kernels:

We can Kernelize SVM's,
ie write $w = \sum\limits_{i=1}^{n} \alpha_i x^{(i)}$

$\alpha_i = 0$ if $x^{(i)}$ is not support vector

$\Rightarrow$ Test time: only evaluate #support vector Kernel calls
instead of $n$

Takeaway: In practice, to use Kernels, use SVM
don't use Kernelized logistic regression