# ① Regularization

Recall logistic regression

Basic/unregularized objective: $L(w) = \sum_{i=1}^{n} -\log \sigma \left(y^{(i)} \cdot w^T x^{(i)}\right)$

<span style="color:blue">margin<br>larger better</span>

<span style="color:orange">Goal: maximize margin on<br>all training examples</span>

<span style="color:green">trade-off</span>

### L2 Regularization
add $\lambda \|w\|^2$ to objective

<span style="color:blue">constant<br>hyperparam</span>

<span style="color:orange">Goal: have smaller<br>entries in $w$</span>

<span style="color:orange">prefer simple<br>function</span>

### L1 Regularization
add $\lambda \|w\|_1$ to objective

$\sum_{j=1}^{d} |w_j|$

<span style="color:orange">Goal: have some entries<br>of $w = 0$</span>

---

## Maximum a Posterior Inference (MAP) <span style="color:green">← provides justification for $L_2$, $L_2$ reg</span>

parameters $w$, dataset $D$

<span style="color:red">view both as random variables (Bayesian)</span>

<span style="color:blue">(In contrast, in MLE $w$ is not a random variable) (Frequencist)</span>

natural goal: <span style="color:green">prior over $w$</span>

maximize $P(w|D) = \dfrac{P(w)\,P(D|w)}{P(D)}$   <span style="color:red">← same as MLE</span>

<span style="color:red">← normalizing constant<br>doesn't depend on $w$</span>

Suppose $P(w)$ is that each $w_j \sim N(0, \sigma^2)$ independently

$P(w) = \prod_{j=1}^{d} \frac{1}{\sigma\sqrt{2\pi}} \cdot \exp\left(-\frac{w_j^2}{\sigma^2}\right)$

$\log P(w) = \sum_{j=1}^{d} \log\left(\underbrace{\frac{1}{\sigma\sqrt{2\pi}}}_{\text{constant}}\right) - \frac{w_j^2}{\sigma^2} = \text{constant} - \sum_{j=1}^{d} \frac{w_j^2}{\sigma^2}$

$= \text{constant} - \frac{1}{\sigma^2}\|w\|^2$

maximizing $\log p(w) \iff$ minimizing $\|w\|^2$

if $p(w)$ Gaussian, equivalent to $L_2$ regularization

---

**Kernals**: a way to modify existing algorithms

- ☐ kernalized linear regression
- ☐ kernalized logistic regression
  ...

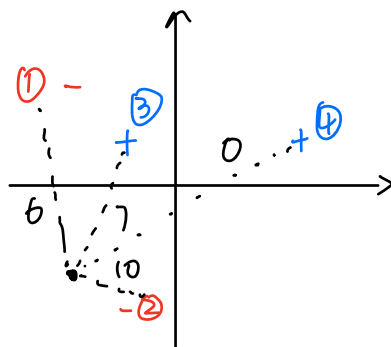Infinite-dimensional features

(by preprocessing data, we can add new features)

| price | Area | Area² | ... |
|-------|------|-------|-----|
|       |      |       |     |

Kernals systematically add many features (possible infinite), but also give a way to work with big feature vectors efficiently.

Make predictions by measuring similarity of test example to each training example. (similar to k-NN)



example weights:

$a_1 = -1$
$a_2 = -1$
$a_3 = 1$
$a_4 = 1$

predict: $-1\cdot 6 -1\cdot 10 +1\cdot 0 +1\cdot 7$
$= -9 \implies -1$

In kernal logistic regression:

prediction $f(x) = \sum_{i=1}^{\Lambda} a_i \, k(x, x^{(i)})$

sum over training examples ⌐ weights learned ⌐ kernal function measures how similar two inputs are

if $f(x) > 0$, predict $1$

if $f(x) < 0$, predict $-1$

Consider $k(x, z) = x^T z$

Captures some notion of similarity, if $x = z$, $k(x, z) > 0$

if point in opposite direction, $k(x, z) < 0$

Logistic regression algorithm can be written in terms of kernels only

(any $x$'s only show up in the kernel function)

## Logistic regression via GD

$w^{(0)} \leftarrow 0 \in \mathbb{R}^d$

for $t = 1, \cdots, T$:

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

$w \leftarrow w^{(T)}$

Given $X$, compute $w^T x$ for prediction

① $w$ is a linear combination of $x^{(i)}$'s

"reparametrize" algorithm to update coefficients of this linear combination

## Kernel logistic regression

$a^{(0)} \leftarrow 0 \in \mathbb{R}^n$    define $w^{(t)} = \sum_{i=1}^{n} a_i^{(t)} \cdot x^{(i)}$

for $t = 1, \cdots, T$:

    for $i = 1, \cdots, n$:

$$a_i^{(t)} \leftarrow a_i^{(t-1)} + \eta \cdot \sigma(-y^{(i)} \cdot \underline{w^{(t-1)T} x^{(i)}}) \cdot y^{(i)}$$

return $a = a^{(T)}$

$$= \left( \sum_{j=1}^{n} a_j^{(t-1)} x^{(j)} \right)^T x^{(i)}$$

$$= \sum_{j=1}^{n} a_j^{(t-1)} x^{(j)T} x^{(i)}$$

$$= \sum_{j=1}^{h} a_j^{(t-1)} K(x^{(j)}, x^{(i)})$$

prediction: $W^T x$ for test input $x$

$$= \sum_{J=1}^{n} a_J K(x^{(J)}, x)$$

If we can compute kernal between any 2 x's, we don't have to store x's themselves.