

$$\hat{w}_{ML} = \underbrace{(X^T X)^{-1} X^T y}_{\text{Random part}}$$

Linear transformations.

$$\begin{aligned} \mathbb{E}[\hat{w}_{ML}] &= \mathbb{E}[(X^T X)^{-1} X^T y] \\ &= (X^T X)^{-1} X^T \mathbb{E}[y]. \end{aligned}$$

$$= \underline{(X^T X)^{-1}} \times \underline{(X^T w)}$$

$$= w$$

$\Rightarrow \hat{w}_{ML}$  is unbiased!

$$y_i = w^T x_i + \epsilon$$

$$y_i = \mathcal{N}(w^T x_i, \sigma^2)$$

$$\mathbb{E}(y_i) = w^T x_i$$

$$\mathbb{E}(y) = X^T w$$

what is variance of  $\hat{w}_{ML}$ ?

---

$$\hat{w}_{ML} = (X^T X)^{-1} X^T y = Ay$$

$$[A = (X^T X)^{-1} X^T]$$

$$\text{Cov}(\hat{w}_{ML}) = \mathbb{E} \left[ (\hat{w}_{ML} - \underbrace{\mathbb{E}[\hat{w}_{ML}]}_{w}) (\hat{w}_{ML} - \mathbb{E}[\hat{w}_{ML}])^T \right]$$

$$= \mathbb{E} \left[ (Ay - w)(Ay - w)^T \right]$$

$$= \mathbb{E} \left[ Ayy^T A^T - Ayw^T - wy^T A^T + ww^T \right]$$

on Simplification [ Please do This ]

$$\text{Cov}(\hat{w}_{ML}) = \sigma^2 (X^T X)^{-1}$$

$$\begin{aligned} \mathbb{E} \left( \underbrace{\|\hat{w}_{ML} - w\|^2}_{\text{MSE}} \right) &= \mathbb{E} \left( \underbrace{(\hat{w}_{ML} - w)^T}_{\alpha} \underbrace{(\hat{w}_{ML} - w)}_{\text{blue}} \right) \\ &= \mathbb{E} \left( \text{trace} \left( \underbrace{(\hat{w}_{ML} - w) (\hat{w}_{ML} - w)^T}_{\text{red underline}} \right) \right) \\ &= \text{trace} \left( \mathbb{E} \left[ \underbrace{(\hat{w}_{ML} - w) (\hat{w}_{ML} - w)^T}_{\text{purple underline}} \right] \right) \end{aligned}$$

$a = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_k \end{bmatrix}$   
 $a^T a = \begin{bmatrix} a_1^2 & a_1 a_2 & a_1 a_3 & \dots & a_1 a_k \\ a_2 a_1 & a_2^2 & a_2 a_3 & \dots & a_2 a_k \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ a_k a_1 & a_k a_2 & a_k a_3 & \dots & a_k^2 \end{bmatrix}$  M.S.E.  
 $\text{trace}(a a^T) = a^T a$   
 $\text{trace}(a^T a) = a^T a$

$$= \text{trace} \left( \text{cov}(\hat{w}_{ML}) \right)$$

$$= \text{trace} \left( \sigma^2 (X X^T)^{-1} \right)$$

$$= \sigma^2 \text{trace} (X X^T)^{-1}$$


---

Let the eigenvalues of  $X X^T$  be  $\{\lambda_1, \dots, \lambda_d\}$

Eigenvalues of  $(X X^T)^{-1} = \left\{ \frac{1}{\lambda_1}, \dots, \frac{1}{\lambda_d} \right\}$

⇒ Mean Squared Error

$$\mathbb{E}(\|\hat{\omega}_{ML} - \omega\|^2) = \underbrace{\sigma^2}_{\downarrow} \underbrace{\left( \sum_{i=1}^d \frac{1}{\lambda_i} \right)}_{\text{}} \underbrace{\quad}_{\text{}}$$

---

$$\hat{\omega}_{ML} = \underbrace{(X^T X)^{-1}}_{\text{}} X y$$

Consider the following modified estimator

$$\hat{\omega}_{\text{modified}} = (X^T X + \lambda I)^{-1} X y$$

- For some matrix  $A$ , say e.val are  $\{\lambda_1, \dots, \lambda_d\}$

$$Au_k = \lambda_k u_k \quad \text{for some } u_k \in \mathbb{R}^d$$

- What are e.values of  $A + \lambda I$

$$\begin{aligned} (A + \lambda I) u_k &= Au_k + \lambda u_k \\ &= \lambda_k u_k + \lambda u_k \\ &= \underline{(\lambda_k + \lambda)} u_k \end{aligned}$$

## EXISTENCE THEOREM (INFORMAL)

$\exists \lambda \in \mathbb{R}_+$  s.t.

$$\hat{w}_{\text{modified}} = (X^T X + \lambda I)^{-1} X^T y \quad \text{has lesser}$$

m.s.e than  $\hat{w}_{ML}$

## Ridge Regression: Biased Estimation for Nonorthogonal Problems

Arthur E. Hoerl & Robert W. Kennard

SO FAR

LINEAR REGRESSION

$$y = X^T w + \epsilon, \quad \epsilon \sim N(0, \sigma^2 I)$$

$$w^* = \underline{(X^T X)^{-1} X^T y}$$

$$\hat{w}_{ML} = \underline{(X^T X)^{-1} X^T y}$$

$$\mathbb{E}[\hat{w}_{ML}] = w$$

$$\text{Cov}(\hat{w}_{ML}) = \sigma^2 (X^T X)^{-1}$$

$$\text{MSE} = \underline{\sigma^2 \text{trace} (X^T X)^{-1}}$$

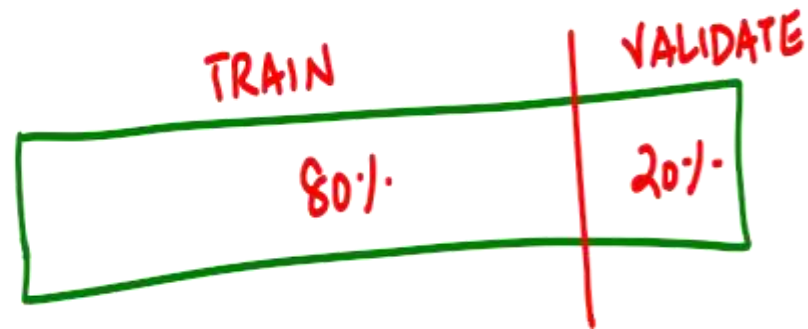
$$\boxed{\hat{w}_{\text{modified}}} = \underline{(X^T X + \lambda I)^{-1} X^T y}$$

$\lambda$  s.t.  $\hat{w}_{\text{mod}}$  has  
smaller m.s.e than  $\hat{w}_{ML}$

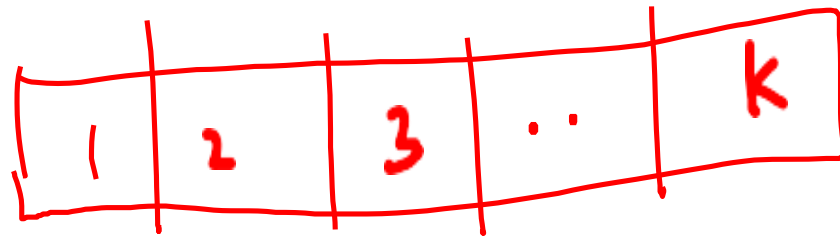
$$\mathbb{E}[\|\hat{w}_{ML} - w\|^2]$$



## CROSS-VALIDATION



K-fold Cross Validation



Leave one out cross validation –  $K = n$

# BAYESIAN MODELING

need

- PRIOR on parameter

- i.e.,  $p(\underline{w})$   $w \in \mathbb{R}^d$

LIKELIHOOD

$$y/x \sim N(w^T x, 1)$$

For simplicity. Can use  $\sigma^2$  as well

PRIOR

$$w \sim \mathcal{N}\left(0, \gamma^2 I\right) \rightarrow \begin{bmatrix} \gamma^2 & & \\ & \gamma^2 & \\ & & \ddots \\ & & & \gamma^2 \end{bmatrix}$$

$\nwarrow \in \mathbb{R}^d$        $\swarrow \in \mathbb{R}^{d \times d}$

As usual

$$P(w | \{(x_1, y_1), \dots, (x_n, y_n)\}) \propto P(\{(x_1, y_1), \dots, (x_n, y_n)\} | w) \cdot P(w)$$

$\uparrow$  posterior       $\uparrow$  Likelihood       $\uparrow$  Prior

$$\propto \left( \prod_{i=1}^n e^{-\frac{(y_i - w^T x_i)^2}{2}} \right) \cdot \left( \prod_{i=1}^d e^{-\frac{w_i^2}{2\gamma^2}} \right)$$

$$\propto \left( e^{-\sum_{i=1}^n \frac{(y_i - w^T x_i)^2}{2}} \right) e^{-\underbrace{\sum_{i=1}^d \frac{w_i^2}{2\gamma^2}}_{\|w\|^2}}$$

$$\log(P(w|\text{Data})) \equiv \bigcirc \sum_{i=1}^n \frac{(y_i - w^T x_i)^2}{2} \bigcirc \frac{1}{2\gamma^2} \|w\|^2$$

## MAP ESTIMATE

$\hat{w}_{\text{MAP}}$

$=$

$\arg \min_w$

$$\sum_{i=1}^n \frac{(y_i - w^T x_i)^2}{2} + \underbrace{\left( \frac{1}{2\gamma^2} \right) \|w\|^2}_{\substack{\text{red arrow: } \frac{1}{2} \|x^T w - y\|^2 \\ \text{blue arrow: } \frac{1}{2\gamma^2} w^T w}}$$

$\rightarrow f(w)$

$$\hat{w}_{\text{MAP}} = \left( X X^T + \frac{1}{\gamma^2} I \right)^{-1} X y$$

CROSS VALIDATE

### CONCLUSION

- MAP estimation for lin. reg with a Gaussian prior  $N(0, \gamma^2 I)$  for  $w$  is "equivalent" to  $\hat{w}_{\text{modified}}$  estimator that we used earlier.

$\hat{w}_R$



$\arg \min_w$

$$\sum_{i=1}^n (w^T x_i - y_i)^2$$



↑  
LOSS

+

$$\lambda \|w\|^2 + \left( \sum_{i=1}^d w_i^2 \right)$$

↑  
REGULARIZATION

HYPER PARAMETER

$$\frac{1}{2\lambda^2}$$

RIDGE  
REGRESSION

## LINEAR REGRESSION

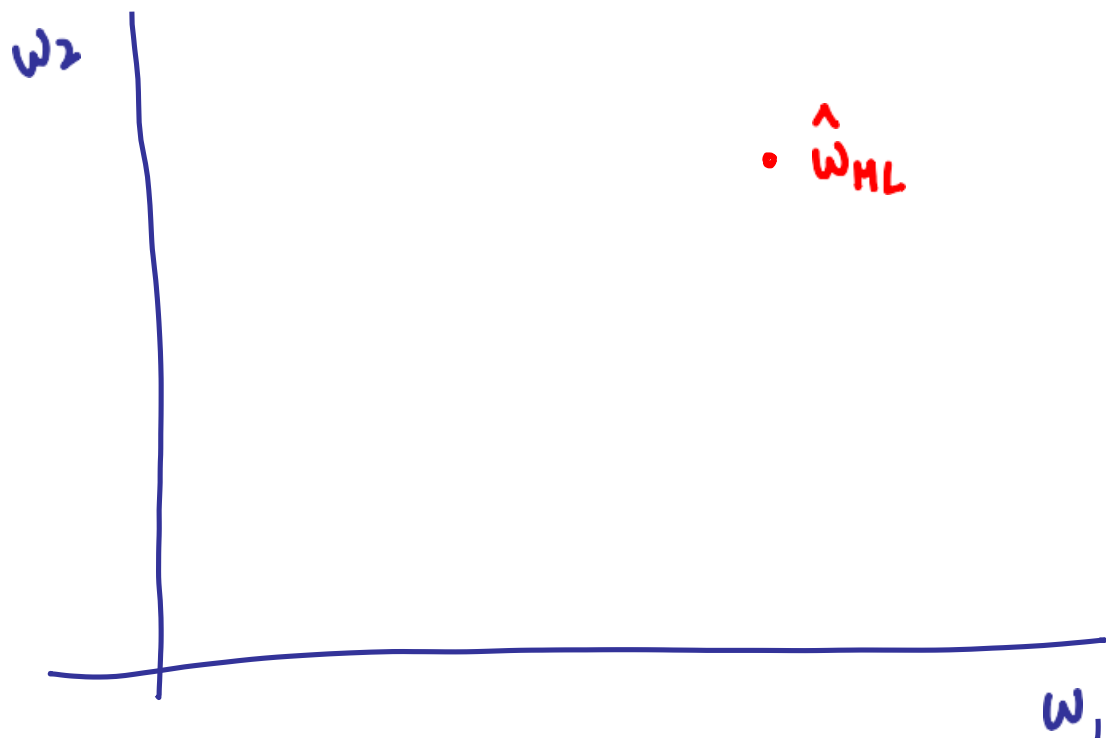
$$\hat{w}_{ML} = \arg \min_{w \in \mathbb{R}^d} \sum_{i=1}^n (w^T x_i - y_i)^2$$

## RIDGE REGRESSION

$$\hat{w}_R = \arg \min_{w \in \mathbb{R}^d} \underbrace{\sum_{i=1}^n (w^T x_i - y_i)^2}_{\text{LOSS}} + \underbrace{\lambda \|w\|^2}_{\text{REGULARIZATION}}$$

HYPER PARAMETER





$$\min_w \sum_{i=1}^n (w^T x_i - y_i)^2 + \lambda \|w\|^2 \rightarrow \textcircled{A}$$

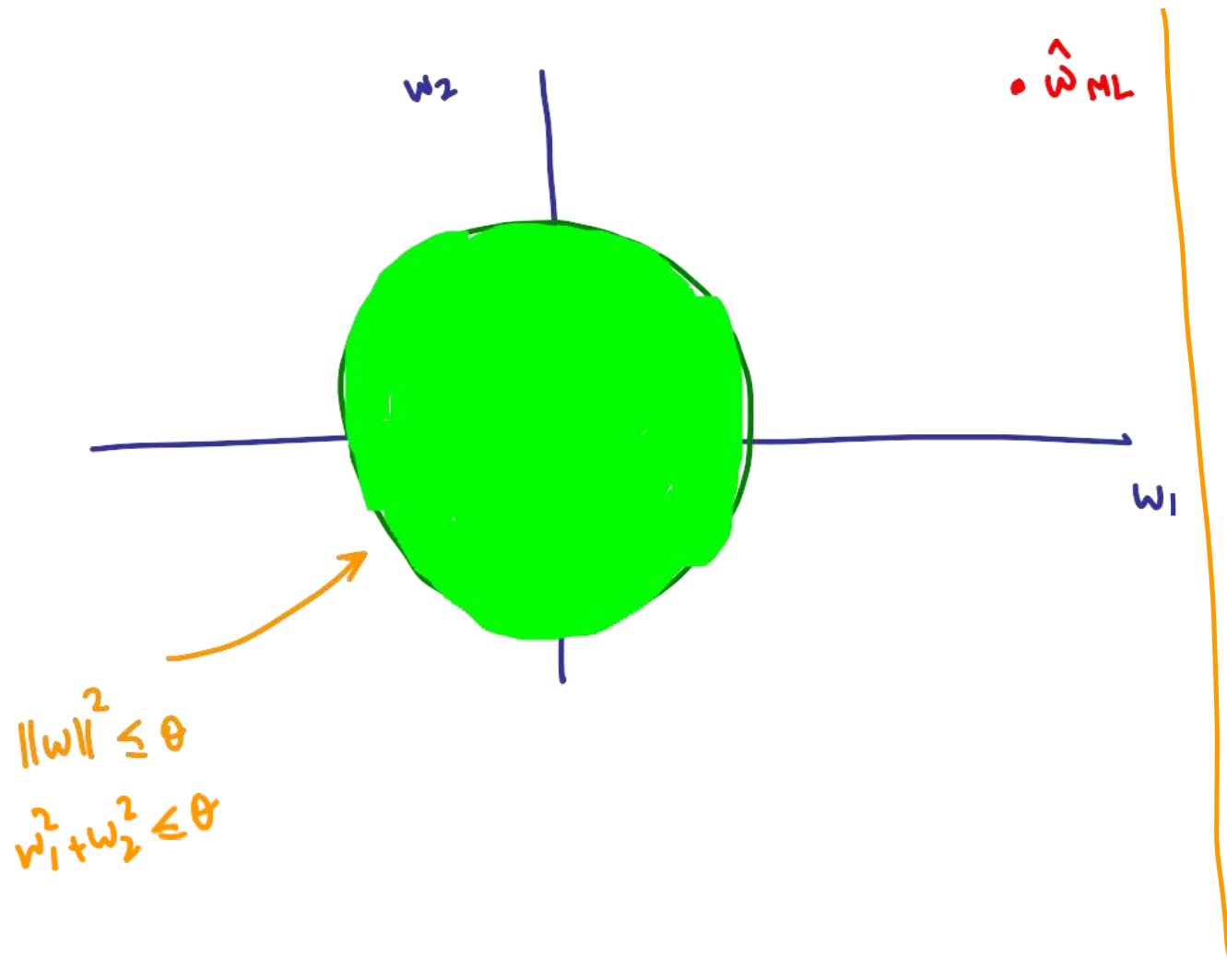
$$= \min_w \sum_{i=1}^n (w^T x_i - y_i)^2 \rightarrow \textcircled{B}$$

$$\text{s.t. } \|w\|^2 \leq \theta$$

↑  
depends on  
 $\lambda$

For every choice of  $\lambda$ ,  $\exists \theta$  (depending on  $\lambda$ ) s.t.

$\textcircled{A}$  and  $\textcircled{B}$  give the same solution.



What is the objective value of linear reg at  $\hat{w}_{ML}$ ?

$$\sum_{i=1}^n \left( \hat{w}_{ML}^T x_i - y_i \right)^2$$

$$= f(\hat{w}_{ML})$$

Consider the following set

$$S_c = \left\{ w \in \mathbb{R}^2 : f(w) = f(\hat{w}_{ML}) + c \right.$$

$c \in \mathbb{R}_+$

$\forall w \in S_c$

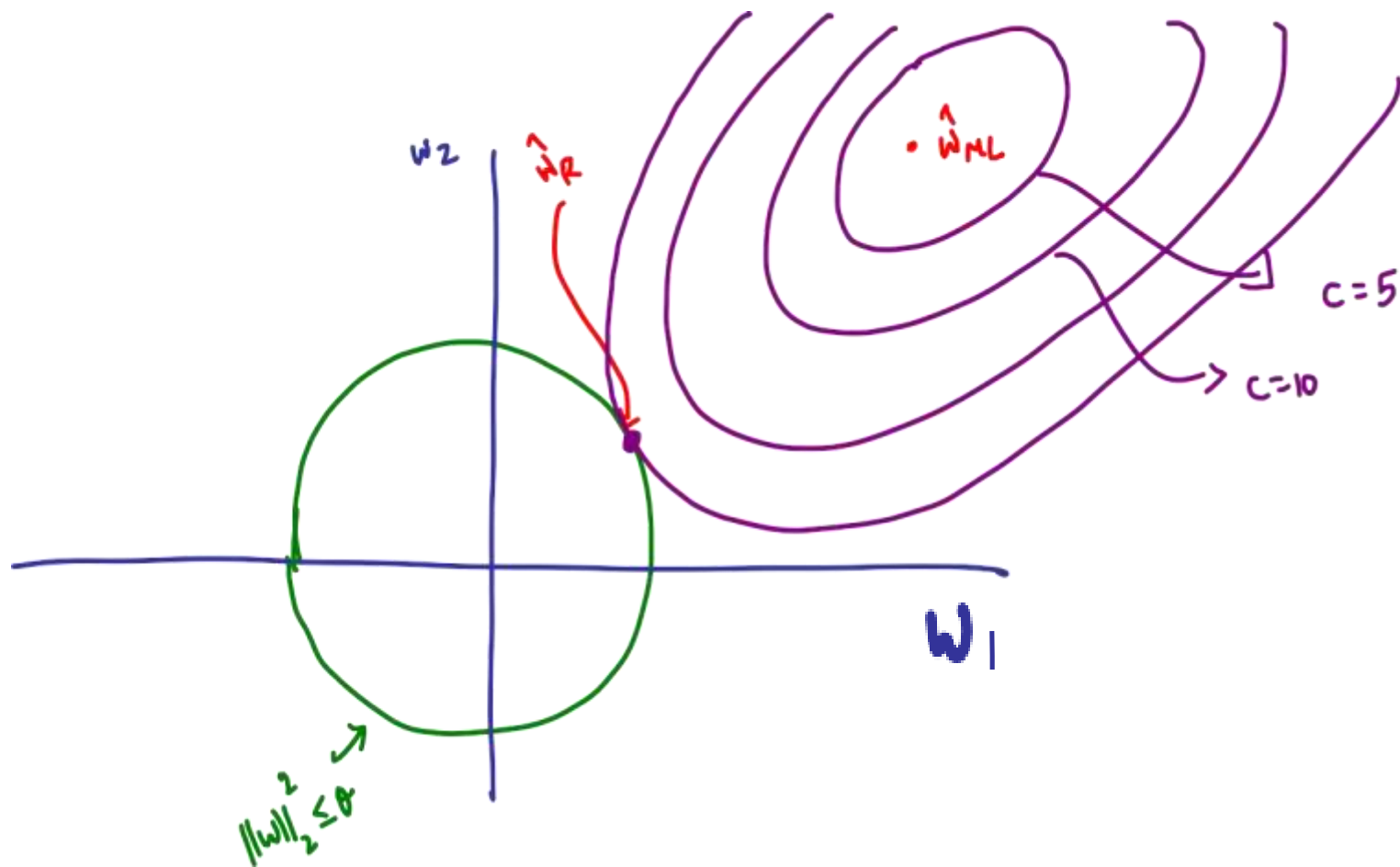
$$\sum_{i=1}^n (w^T x_i - y_i)^2 = \sum_{i=1}^n (\hat{w}_{ML}^T x_i - y_i)^2 + c$$

$$\therefore \underbrace{\|X^T w - y\|^2}_{\text{}} = \|X^T \hat{w}_{ML} - y\|^2 + c$$

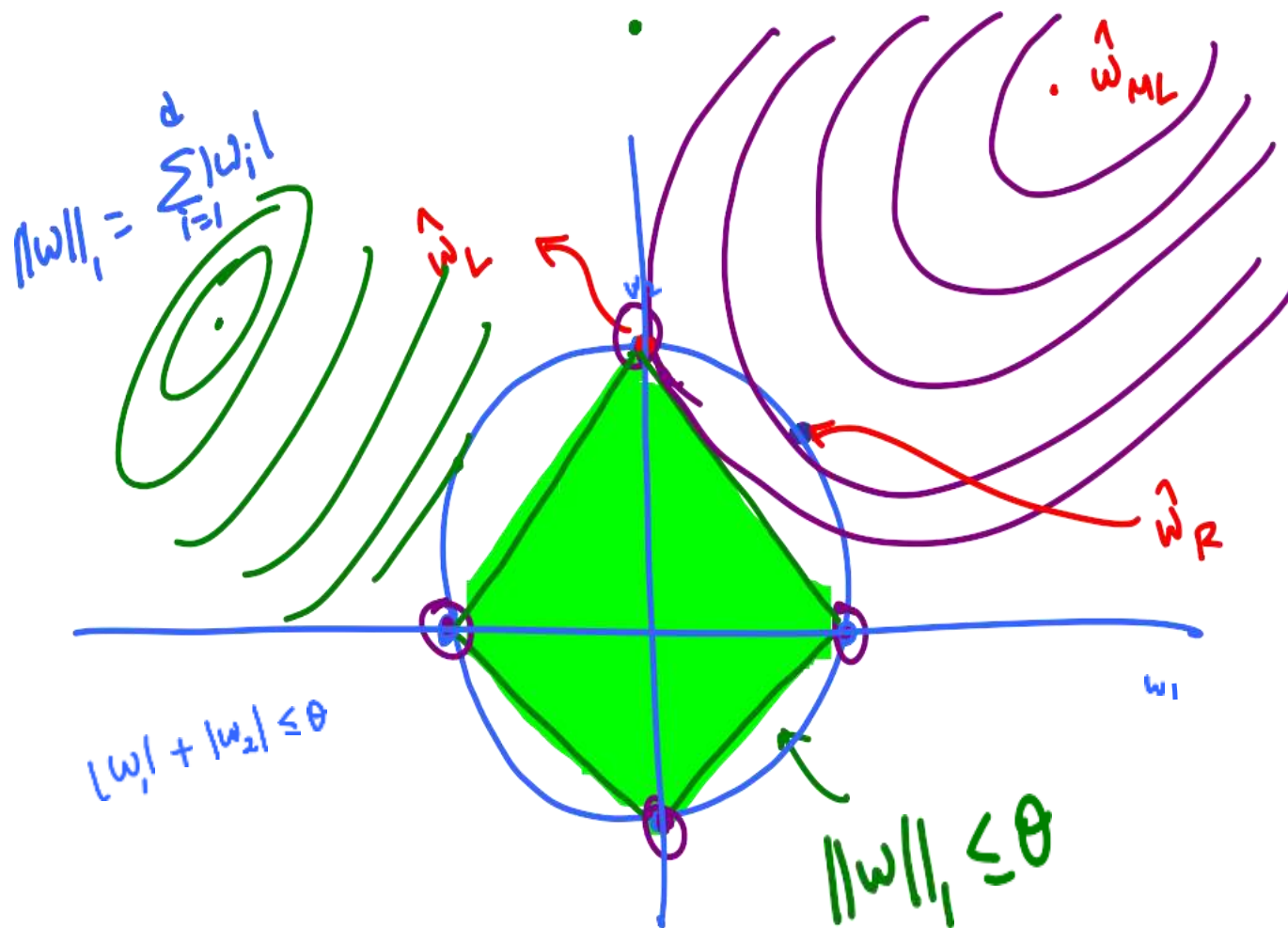
on Simplification [do this]

$$\Rightarrow \underline{(w - \hat{w}_{ML})^T (xx^T) (w - \hat{w}_{ML})} = c' \swarrow$$

depends on  
 $c, (xx^T), \hat{w}_{ML}$   
but does not  
depend on  $w$



- Ridge regression pushes feature weights towards 0 but does not necessarily make it 0.



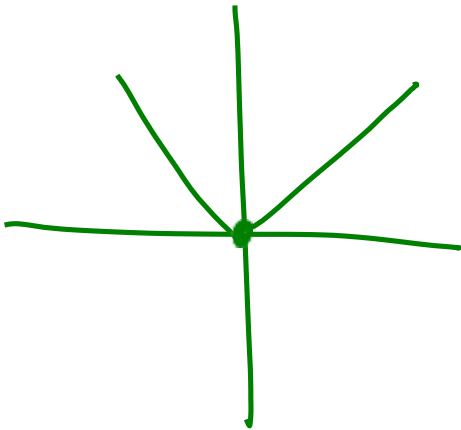
$$\begin{aligned} \min_{w \in \mathbb{R}^d} \quad & \sum_{i=1}^n (w^T x_i - y_i)^2 \\ \text{s.t.} \quad & \|w\|_1 \leq \theta \end{aligned}$$

LASSO

min  
 $w \in \mathbb{R}^d$

$$\sum_{i=1}^n (w^T x_i - y_i)^2 + \lambda \|w\|_1,$$

→ LEAST ABSOLUTE SHRINKAGE &  
SELECTION OPERATOR



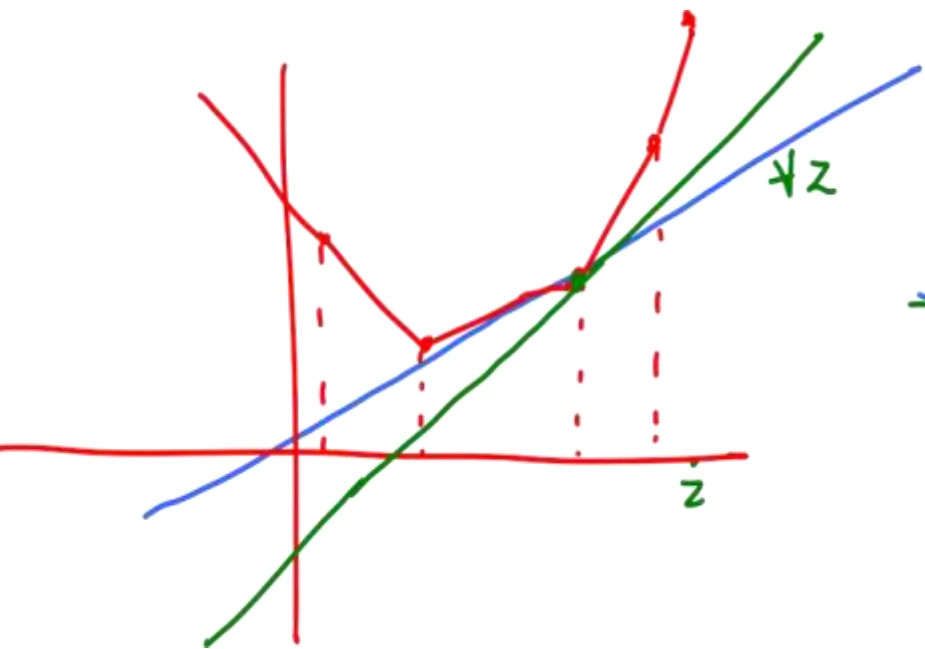
• How to solve this?

→ Sub-gradient descent

## Sub-gradient

A vector  $g \in \mathbb{R}^d$  is a sub-grad of  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  at  $x$  if

$$f(z) \geq \underline{f(x)} + g^T(z-x)$$



If  $f$  is convex, Sub-grad descent converges!



## MODELING NON LINEAR RELATIONSHIPS

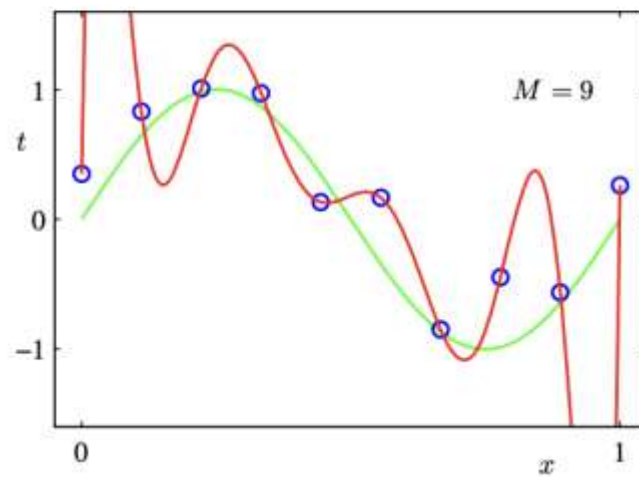
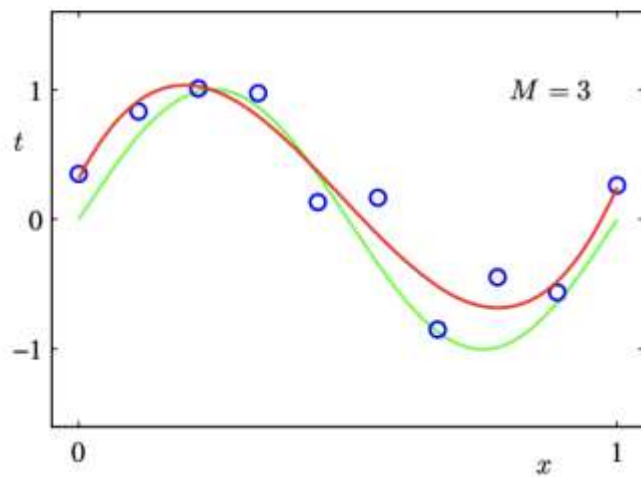
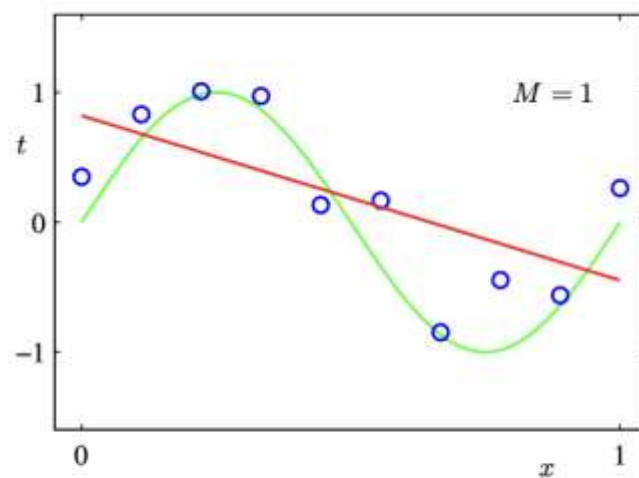
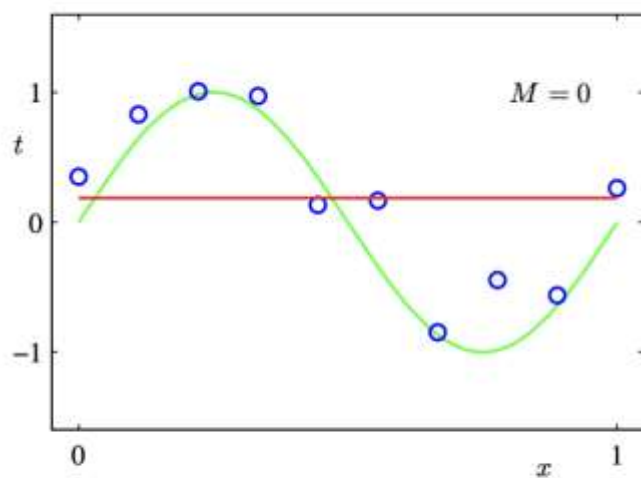
Key idea: Map features to higher dimension.

$$x = [x_1, x_2] \rightarrow \phi(x) = [1, x_1, x_2, x_1^2, x_2^2, x_1x_2]$$

In the above example, learning a linear function in the higher dimension is equivalent to learning a quadratic function in the lower dimension.

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

$$y(x, \mathbf{w}) = w_0 + w_1x + w_2x^2 + \dots + w_Mx^M = \sum_{j=0}^M w_jx^j$$



## Coefficient values increases as M increases!

|         | $M = 0$ | $M = 1$ | $M = 6$ | $M = 9$     |
|---------|---------|---------|---------|-------------|
| $w_0^*$ | 0.19    | 0.82    | 0.31    | 0.35        |
| $w_1^*$ |         | -1.27   | 7.99    | 232.37      |
| $w_2^*$ |         |         | -25.43  | -5321.83    |
| $w_3^*$ |         |         | 17.37   | 48568.31    |
| $w_4^*$ |         |         |         | -231639.30  |
| $w_5^*$ |         |         |         | 640042.26   |
| $w_6^*$ |         |         |         | -1061800.52 |
| $w_7^*$ |         |         |         | 1042400.18  |
| $w_8^*$ |         |         |         | -557682.99  |
| $w_9^*$ |         |         |         | 125201.43   |

## Regularization helps control the co-efficient values

|         | $\ln \lambda = -\infty$ | $\ln \lambda = -18$ | $\ln \lambda = 0$ |
|---------|-------------------------|---------------------|-------------------|
| $w_0^*$ | 0.35                    | 0.35                | 0.13              |
| $w_1^*$ | 232.37                  | 4.74                | -0.05             |
| $w_2^*$ | -5321.83                | -0.77               | -0.06             |
| $w_3^*$ | 48568.31                | -31.97              | -0.05             |
| $w_4^*$ | -231639.30              | -3.89               | -0.03             |
| $w_5^*$ | 640042.26               | 55.28               | -0.02             |
| $w_6^*$ | -1061800.52             | 41.32               | -0.01             |
| $w_7^*$ | 1042400.18              | -45.95              | -0.00             |
| $w_8^*$ | -557682.99              | -91.53              | 0.00              |
| $w_9^*$ | 125201.43               | 72.68               | 0.01              |

$$x = [f_1 \quad f_2 \quad f_3 \quad f_4]$$

$\Downarrow$  Cubic relation

$$\phi(x) = \left[ 1 \quad \underbrace{f_1 \quad f_2 \quad f_3 \quad f_4}_{1 + 4} \quad + \quad 4c^2 + 4c^3 \right]$$

$$x \in d$$

$$\phi(x) \in d(d^p)$$

ISSUE :

$\phi(x) \in \mathbb{R}^D \rightarrow$  might be too large

### EXAMPLE

$$x = [f_1 \quad f_2]$$

$$x' = [g_1 \quad g_2]$$

$$\begin{aligned} \boxed{\left( \underline{x x' + 1} \right)^2} &= \left( [f_1 \quad f_2] \begin{bmatrix} g_1 \\ g_2 \end{bmatrix} + 1 \right)^2 \\ &= \left( f_1 g_1 + f_2 g_2 + 1 \right)^2 \\ &= \underline{f_1^2 g_1^2} + \underline{f_2^2 g_2^2} + 1 + 2 f_1 g_1 f_2 g_2 \\ &\quad + \underline{2 f_1 g_1} + 2 f_2 g_2 \end{aligned}$$

$$\left( \underline{\underline{\tau x' + 1}} \right)^2$$

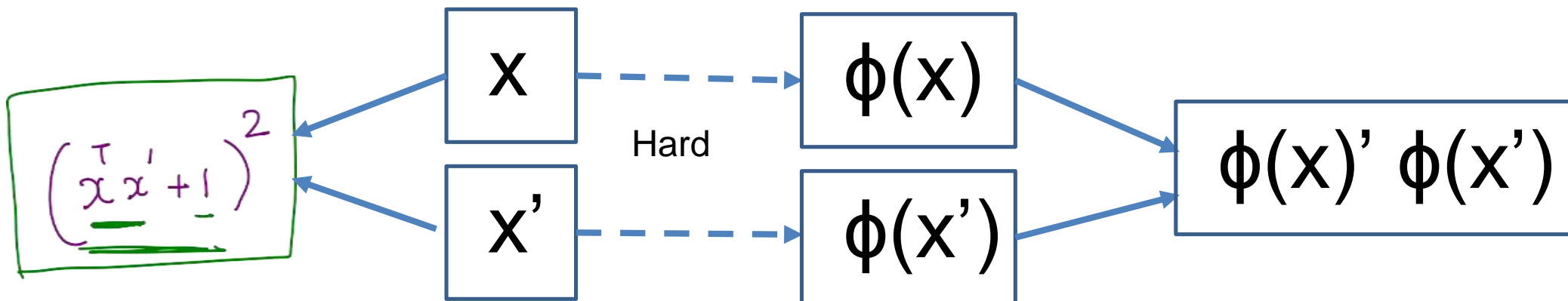
$$= \left[ \underline{f_1}^2, \underline{f_2}^2 \right]$$

$$\uparrow \phi(x)^\tau \phi(x')$$

$$\sqrt{2} f_1, f_2$$

$$\underline{\sqrt{2} f_1}$$

$$\sqrt{2} f_2 \begin{bmatrix} g_1^2 \\ g_2^2 \\ 1 \\ \sqrt{2} g_1 g_2 \\ \sqrt{2} g_1 \\ \sqrt{2} g_2 \end{bmatrix}$$



## Summary/Questions

- > To capture non-linear relationships, one can “create” non-linear functions of features.
- > But the number of features to create grows exponentially with the degree of non-linearity  $p$  that we wish to capture ( $d^p$ )
- > For  $d = 2$  and  $p = 2$ , it appears there is a trick to get around this.
- > Is this trick general enough to be useful the general case as well? (i.e., for any  $d$  and any  $p$ ?)

## MORE EXAMPLES

Polynomial map

$$k(x, x') = \left( \underline{x^T x' + 1} \right)^p$$

for some  $p \geq 1$

→ Can be shown to be a "valid" function.

i.e.,  $\exists \phi: \mathbb{R}^d \rightarrow \mathbb{R}^D$  such that

$$k(x, x') = \phi(x)^T \phi(x')$$

EXERCISE

Compute  $\phi$   
for  $p=3$ ,  
 $p=4$ .



## MORE EXAMPLES

$$K(x, x') = \exp\left(-\frac{\|x - x'\|^2}{2\sigma^2}\right)$$

for some  $\sigma > 0$

RADIAL BASIS FUNCTION

→ Can also shown to be a "valid" map

→ Interestingly,  $\phi$  in this case maps  $x$  to an "infinite" dimensional space

[ Technicalities aside  
think of  $\phi(x)$  as a  
function and dot product as  
integrals ]

## KERNEL FUNCTION

Any function  $K: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$  which is a "valid" map is a kernel function

$$K(x, x') = (x^T x' + 1)^p \rightarrow \text{Polynomial kernel}$$

$$K(x, x') = \exp\left(-\frac{\|x - x'\|^2}{2\sigma^2}\right) \rightarrow \text{Gaussian kernel / Radial basis kernel}$$

Question:

Given a function

$\mathbb{R} \cdot \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ , how can we say it's a

valid kernel?

METHOD 1.

Explicitly construct a  $\phi$  map

[might be hard sometimes]

## METHOD 2: MERCER'S THEOREM

A function  $k: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$  is a valid kernel

if and only if

(a)  $k$  is symmetric.

(b) For any dataset  $\{x_1, \dots, x_n\}$ , the

matrix  $K \in \mathbb{R}^{n \times n}$

$$K_{ij} = k(x_i, x_j)$$

is

POSITIVE  
SEMI  
DEFINITE

eigenvalues  
of  
 $K$  are all  
non-negative.

**How does kernel help in solving the non-linear regression problem?**

$$\begin{array}{c} x_1 \dots x_n \\ y_1 \dots y_n \end{array}$$


---

$$\rightarrow \hat{\underline{w}}_{ML}$$

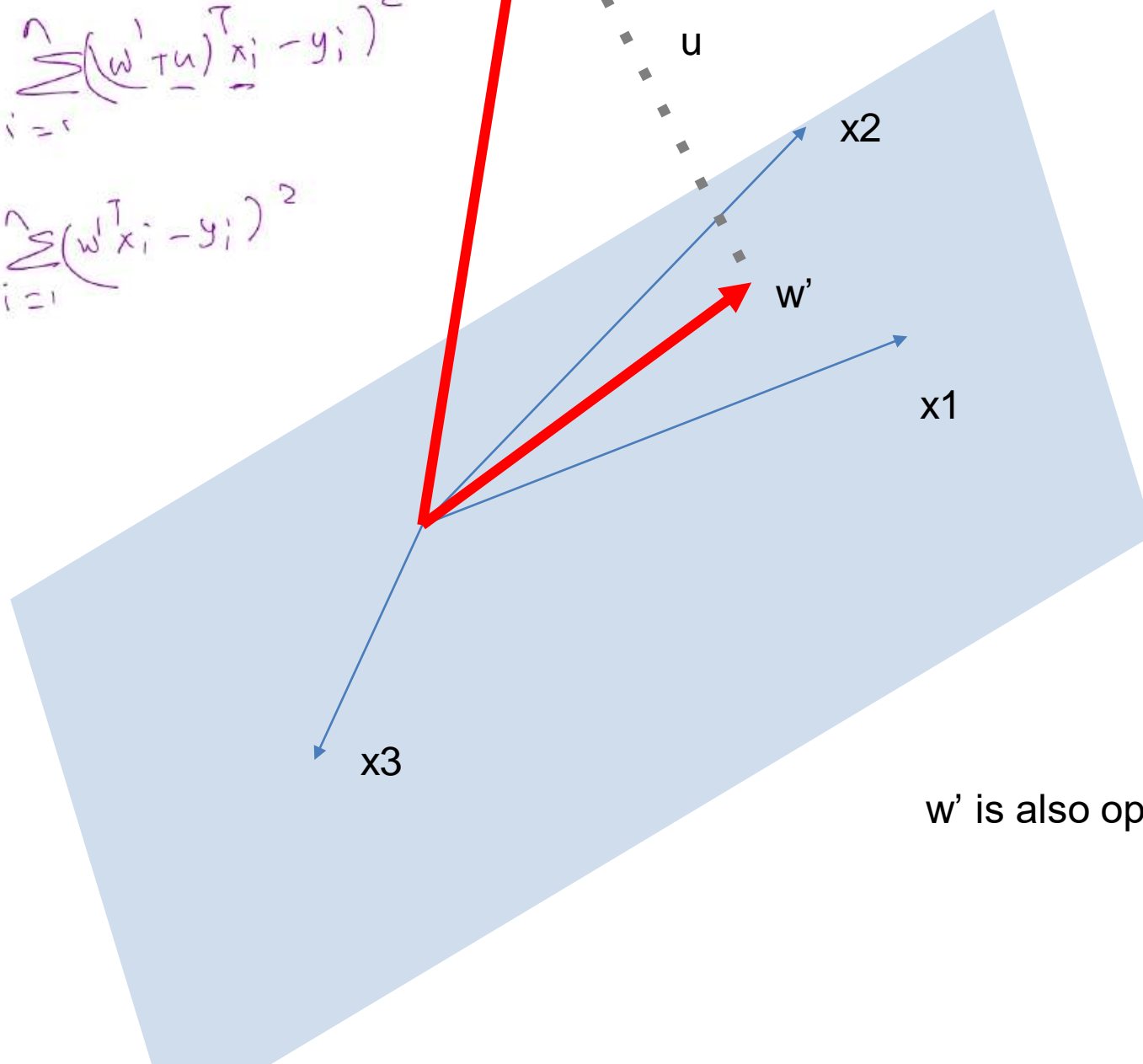
$$\hat{y} = \hat{w}_{ML}^T t$$

$$\text{test point} \\ t \in \mathbb{R}^d$$

Claimed optimal  $w$

$$w = w' + u$$

$$\begin{aligned}\sum_{i=1}^n (w^T x_i - y_i)^2 &= \sum_{i=1}^n (w' + u)^T x_i - y_i)^2 \\ &= \sum_{i=1}^n (w'^T x_i - y_i)^2\end{aligned}$$



$$\hat{w}_{ML} = (XX^T)^{-1}Xy$$

$$\hat{w}_{ML} = X\alpha$$

$$X\alpha = (XX^T)^{-1}Xy$$

$$X^TXX^TX\alpha = X^TXX^T(XX^T)^{-1}Xy$$

$$(X^TX)^2\alpha = X^TXy$$

$$\alpha = (X^TX)^{-1}y$$

**KERNEL**

kernel-regression  $\leftarrow$

given a kernel  $K \in \mathbb{R}^{n \times n}$ ,

obtain  $\underline{\alpha} = \underline{K}^{-1}y \leftarrow \in \mathbb{R}^n$

For a new point  $t$

$$\hat{y} = \sum_{i=1}^n \alpha_i \underbrace{K(t, x_i)}$$