

Training Issues

- (1. Overfitting Bias_Var)
- (2. Optimization)

Bias-Var

```
In [1]: import numpy as np
import tensorflow as tf
import matplotlib.pyplot as plt
%matplotlib inline

In [1]: # Data for Training
n_data = 123
x_data = np.random.normal(1.0, 2.0, n_data)
y_data = x_data * 0.45 + 0.6 + np.random.normal(0.0, 0.1, n_data)
plt.scatter(x_data, y_data)
plt.show()

In [1]: # Model
w = tf.Variable(tf.random_uniform([-1.0, 1.0]))
b = tf.Variable(tf.zeros([1]))
x = tf.placeholder(tf.float32, n_data) # check by "tab" or "shift+tab"
y = w * x + b

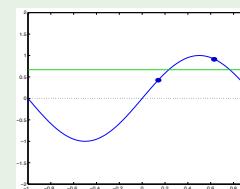
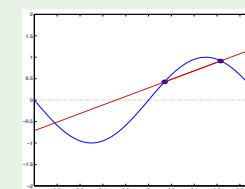
In [1]: # Error function (loss/cost ...)
loss = tf.reduce_mean(tf.square(y - y_data))

sess = tf.Session()
sess.run(init)

In [1]: # Optimization Scheme (or training scheme)
train = tf.train.GradientDescentOptimizer(0.01).minimize(loss)
init = tf.global_variables_initializer()

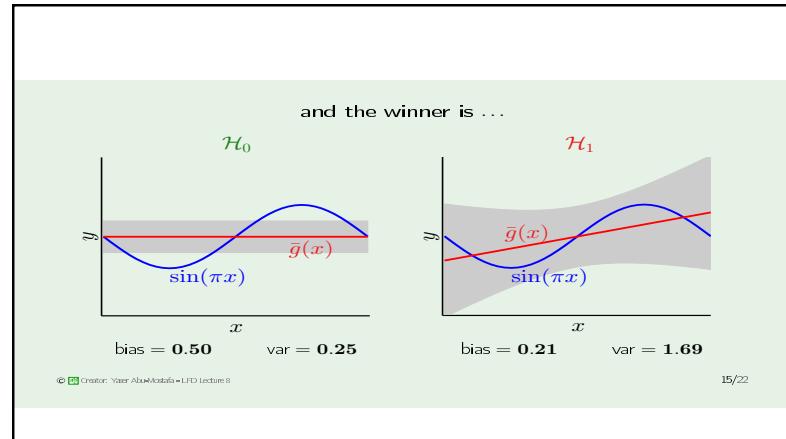
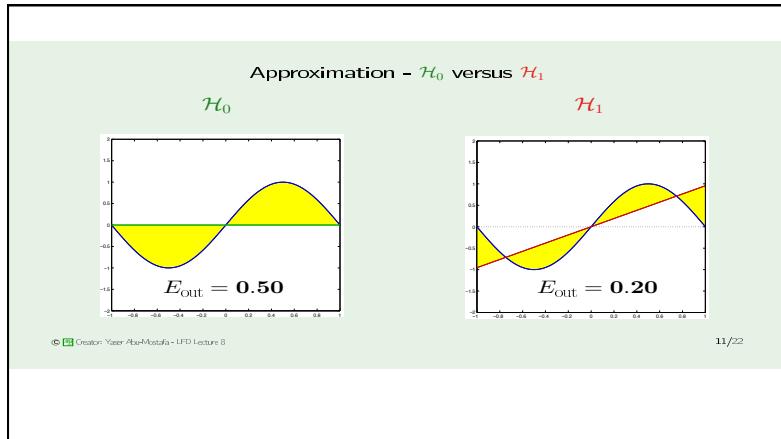
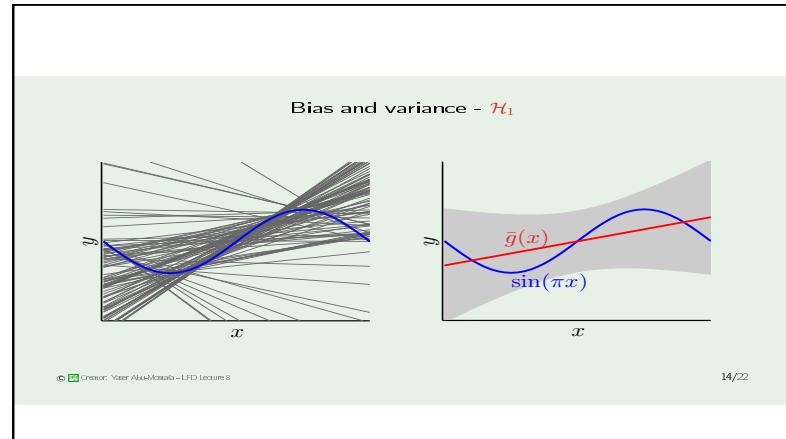
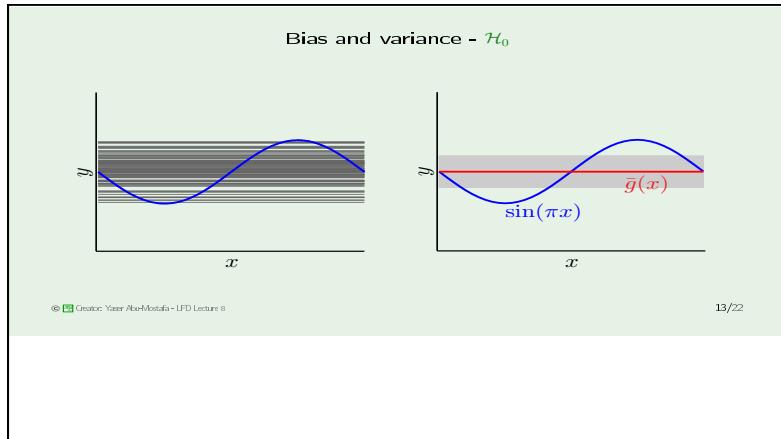
In [1]: # Training
sess = tf.Session()
sess.run(init)
for epoch in range(200):
    sess.run(train, feed_dict={x:x_data})
    if epoch % 50 == 0:
        print(epoch, sess.run(w), sess.run(b))

In [1]: # Evaluation (Checking the ML results)
plt.plot(x_data, y_data, 'ro')
plt.plot(x_data, sess.run(w) * x_data + sess.run(b))
plt.show()
```

Learning - \mathcal{H}_0 versus \mathcal{H}_1 \mathcal{H}_0  \mathcal{H}_1 

© Creator: Yaser Alaa-Katibah-LFD lecture 8

12/22



Bias and variance

$$\mathbb{E}_{\mathcal{D}} \left[(g^{(\mathcal{D})}(\mathbf{x}) - f(\mathbf{x}))^2 \right] = \underbrace{\mathbb{E}_{\mathcal{D}} \left[(g^{(\mathcal{D})}(\mathbf{x}) - \bar{g}(\mathbf{x}))^2 \right]}_{\text{var}(\mathbf{x})} + \underbrace{\left(\bar{g}(\mathbf{x}) - f(\mathbf{x}) \right)^2}_{\text{bias}(\mathbf{x})}$$

Therefore, $\mathbb{E}_{\mathcal{D}} [E_{\text{out}}(g^{(\mathcal{D})})] = \mathbb{E}_{\mathbf{x}} \left[\mathbb{E}_{\mathcal{D}} \left[(g^{(\mathcal{D})}(\mathbf{x}) - f(\mathbf{x}))^2 \right] \right]$

$$= \mathbb{E}_{\mathbf{x}} [\text{bias}(\mathbf{x}) + \text{var}(\mathbf{x})]$$

$$= \text{bias} + \text{var}$$

© Creator: Yaser Abu-Mostafa - LFD Lecture 8 8/22

The curves

Simple Model

Complex Model

© Creator: Yaser Abu-Mostafa - LFD Lecture 8 19/22

VC versus bias-variance

VC analysis

bias-variance

© Creator: Yaser Abu-Mostafa - LFD Lecture 8 20/22

Linear regression case

Noisy target $y = \mathbf{w}^* \mathbf{x} + \text{noise}$

Data set $\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$

Linear regression solution: $\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$

In-sample error vector = $\mathbf{X}\mathbf{w} - \mathbf{y}$

'Out-of-sample' error vector = $\mathbf{X}\mathbf{w} - \mathbf{y}'$

© Creator: Yaser Abu-Mostafa - LFD Lecture 8 21/22

Linear regression case

Noisy target $y = \mathbf{w}^* \mathbf{x} + \text{noise}$

Data set $\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$

Linear regression solution: $\mathbf{w} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$

In-sample error vector = $\mathbf{X}\mathbf{w} - \mathbf{y}$

'Out-of-sample' error vector = $\mathbf{X}\mathbf{w} - \mathbf{y}'$

© [Creative Commons BY-SA] Yaser Abu-Mostafa - LFD Lecture 8

21/22

Learning curves for linear regression

Best approximation error = σ^2

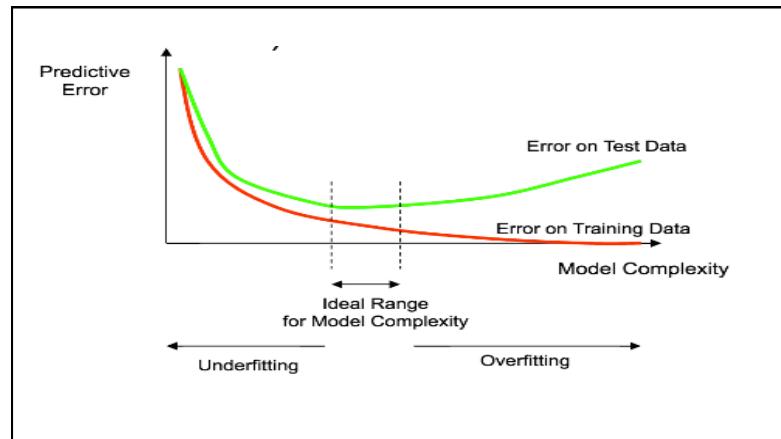
Expected in-sample error = $\sigma^2 \left(1 - \frac{d+1}{N}\right)$

Expected out-of-sample error = $\sigma^2 \left(1 + \frac{d+1}{N}\right)$

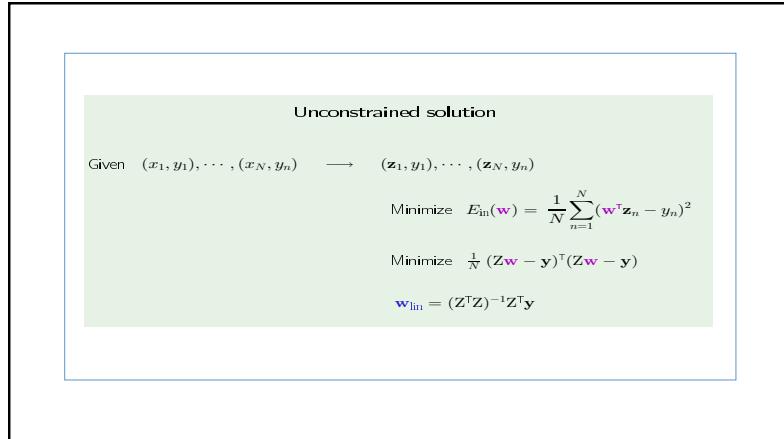
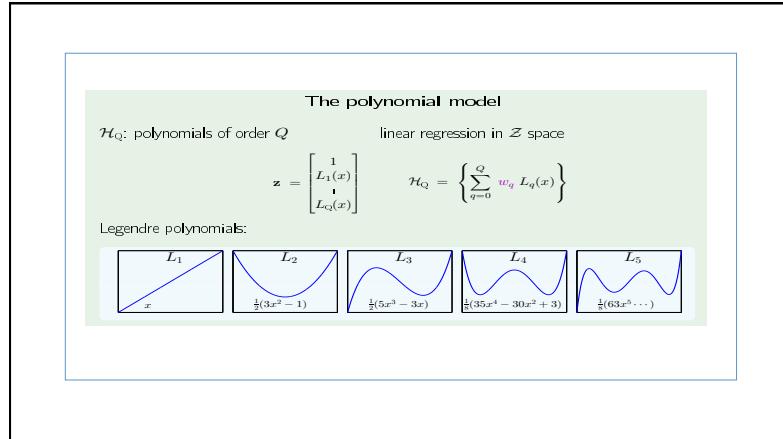
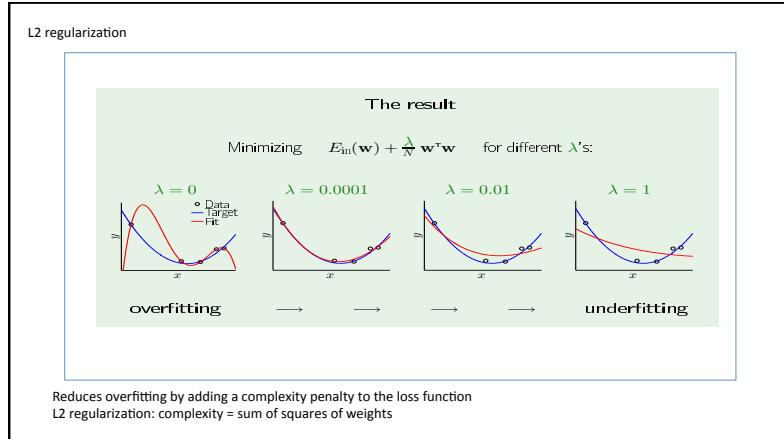
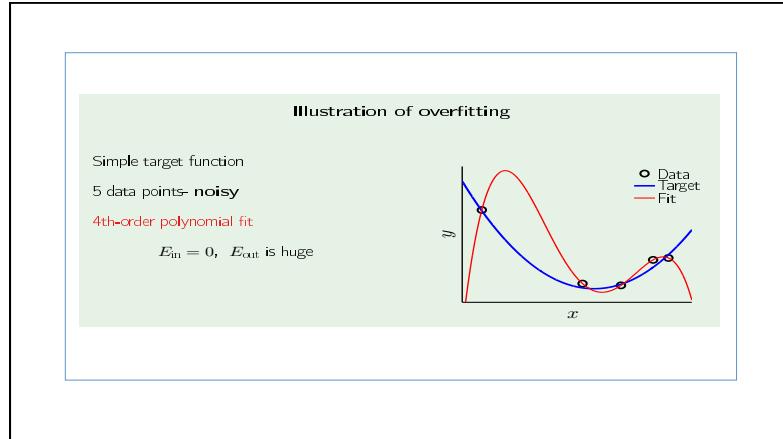
Expected generalization error = $2\sigma^2 \left(\frac{d+1}{N}\right)$

© [Creative Commons BY-SA] Yaser Abu-Mostafa - LFD Lecture 8

22/22



Regularization



Augmented error

Minimizing $E_{\text{aug}}(\mathbf{w}) = E_{\text{in}}(\mathbf{w}) + \frac{\lambda}{N} \mathbf{w}^T \mathbf{w}$

$$= \frac{1}{N} (\mathbf{Z}\mathbf{w} - \mathbf{y})^T (\mathbf{Z}\mathbf{w} - \mathbf{y}) + \frac{\lambda}{N} \mathbf{w}^T \mathbf{w} \quad \text{unconditionally}$$

— solves —

Minimizing $E_{\text{in}}(\mathbf{w}) = \frac{1}{N} (\mathbf{Z}\mathbf{w} - \mathbf{y})^T (\mathbf{Z}\mathbf{w} - \mathbf{y})$

subject to: $\mathbf{w}^T \mathbf{w} \leq C$ \longleftarrow VC formulation

Weight ‘decay’

Minimizing $E_{\text{in}}(\mathbf{w}) + \frac{\lambda}{N} \mathbf{w}^T \mathbf{w}$ is called weight decay. Why?

Gradient descent:

$$\mathbf{w}(t+1) = \mathbf{w}(t) - \eta \nabla E_{\text{in}}(\mathbf{w}(t)) - 2 \eta \frac{\lambda}{N} \mathbf{w}(t)$$

$$= \mathbf{w}(t) (1 - 2\eta \frac{\lambda}{N}) - \eta \nabla E_{\text{in}}(\mathbf{w}(t))$$

Applies in neural networks:

$$\mathbf{w}^T \mathbf{w} = \sum_{l=1}^L \sum_{i=0}^{d^{(l-1)}} \sum_{j=1}^{d^{(l)}} \left(w_{ij}^{(l)} \right)^2$$

Variations of weight decay

Emphasis of certain weights: $\sum_{q=0}^Q \gamma_q w_q^2$

Examples: $\gamma_q = 2^q \implies$ low-order fit
 $\gamma_q = 2^{-q} \implies$ high-order fit

Neural networks: different layers get different γ 's

Tikhonov regularizer: $\mathbf{w}^T \Gamma^T \Gamma \mathbf{w}$

Andrew Ng

Debugging a learning algorithm:

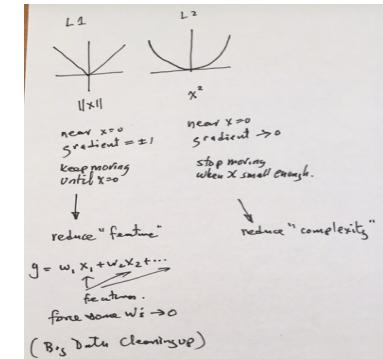
Suppose you have implemented regularized linear regression to predict housing prices. However, when you test your hypothesis in a new set of houses, you find that it makes unacceptably large errors in its prediction. What should you try next?

- Get more training examples \rightarrow fixes high variance
- Try smaller sets of features \rightarrow fixes high variance
- Try getting additional features \rightarrow fixes high bias
- Try adding polynomial features ($x_1^2, x_2^2, x_1 x_2, \text{etc}$) \rightarrow fixes high bias.
- Try decreasing $\lambda \rightarrow$ fixes high bias
- Try increasing $\lambda \rightarrow$ fixes high variance

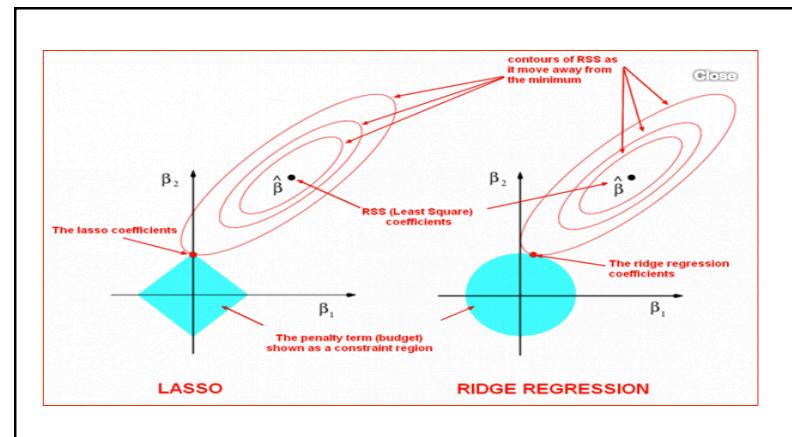
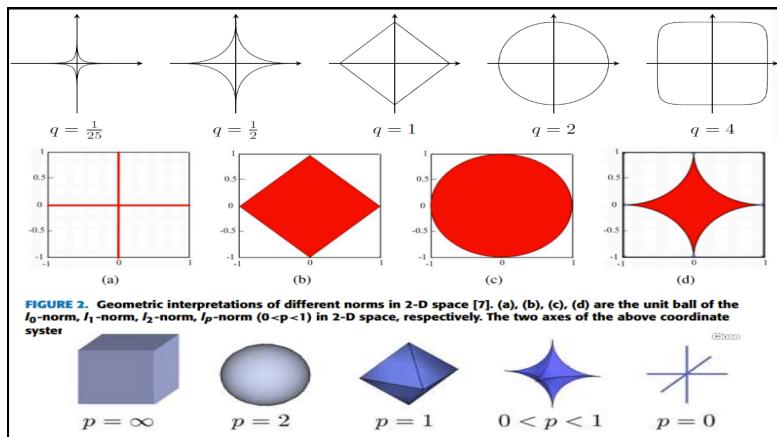
Regularization

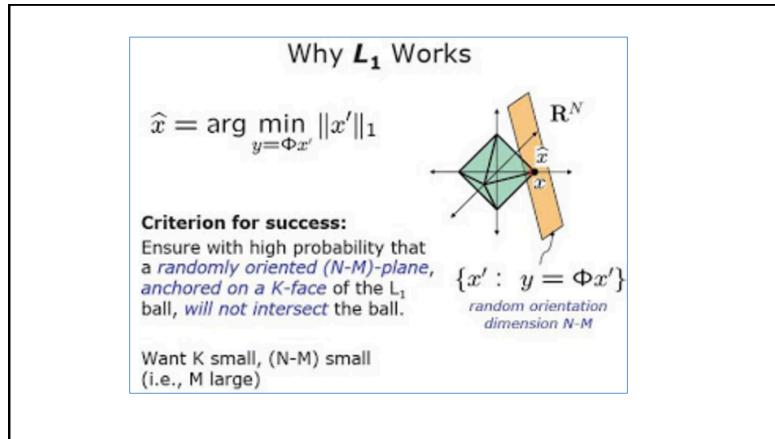
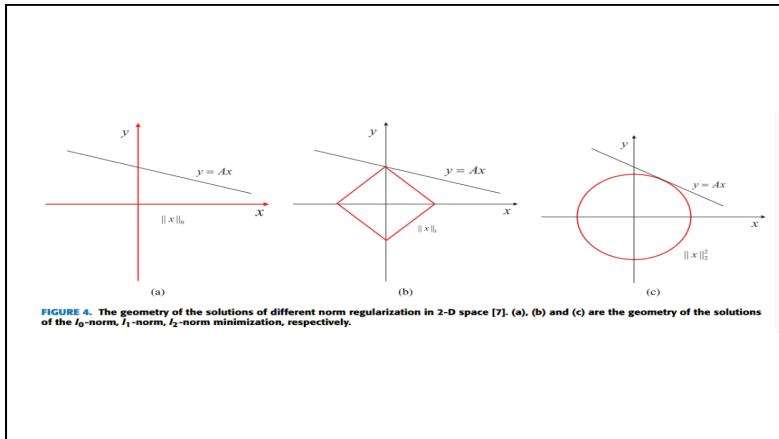
(sparsity)

L1 (LASSO) vs L2 regularization



L1 (LASSO=least absolute shrinkage and selection operator) leads to SPARSE parameterization.





Regularization

(kernel view – kernel linear regression)

Linear regression vs. Kernel regression

The Basic Idea

feature mapping $\phi: R^d \rightarrow R^D$
 $(x_i, z_i) \mapsto (z_i, z_i^2, \dots, z_i^D)$

original space X feature space H

$\hat{y}_j = \omega^\top x_j$

$\hat{y}_j = \omega^\top z_j$

$\hat{y}_j = \omega^\top \phi(x_j)$

■ Linear Kernel

$$\kappa(x, z) = \langle x, z \rangle$$

■ Polynomial Kernel

$$\kappa(x, z) = (\langle x, z \rangle + 1)^r, \quad r \in \mathbb{Z}^+$$

■ RBF (Gaussian) Kernel

$$\kappa(x, z) = \exp\left(\frac{-\|x - z\|^2}{2\sigma^2}\right), \quad \sigma \in \mathbb{R} - \{0\}$$

Regularization

(kernel view -

Classical algorithm:
Regularization in RKHS (eg. kernel machines)

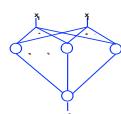
$$\min_{f \in H} \left[\frac{1}{n} \sum_{i=1}^n V(f(x_i) - y_i) + \lambda \|f\|_K^2 \right]$$

implies

$$f(\mathbf{x}) = \sum_i \alpha_i K(\mathbf{x}, \mathbf{x}_i)$$

Remark (for later use):

Classical kernel machines correspond to shallow networks



Proposition 8. Let $\mathbf{z} \in \mathbb{Z}^m$ and $\gamma \in \mathbb{R}$, $\gamma > 0$. The empirical target, i.e. the function $f_{\gamma, \mathbf{z}} = f_{\mathbf{z}}$ minimizing the regularized empirical error

$$\frac{1}{m} \sum_{i=1}^m (y_i - f(x_i))^2 + \gamma \|f\|_K^2$$

over $f \in \mathcal{H}_K$, may be expressed as

$$f_{\mathbf{z}}(x) = \sum_{i=1}^m a_i K(x, x_i)$$

where $a = (a_1, \dots, a_m)$ is the unique solution of the well-posed linear system in \mathbb{R}^m

$$(\gamma m \text{Id} + K[\mathbf{x}])a = \mathbf{y}.$$

Here, we recall, $K[\mathbf{x}]$ is the $m \times m$ matrix whose (i, j) entry is $K(x_i, x_j)$, $\mathbf{x} = (x_1, \dots, x_m) \in X^m$, and $\mathbf{y} = (y_1, \dots, y_m) \in Y^m$ such that $\mathbf{z} = ((x_1, y_1), \dots, (x_m, y_m))$.

Proof. Let $H(f) = \frac{1}{m} \sum_{i=1}^m (y_i - f(x_i))^2 + \gamma \|f\|_K^2$ and write, for any $f \in \mathcal{H}_K$, $f = \sum_{k=1}^{\infty} c_k \phi_k$. Recall that $\|f\|_K^2 = \sum_{k=1}^{\infty} \frac{c_k^2}{\lambda_k}$.

For every $k \geq 1$, $\frac{\partial H}{\partial c_k} = \frac{1}{m} \sum_{i=1}^m -2(y_i - f(x_i))\phi_k(x_i) + 2\gamma \frac{c_k}{\lambda_k}$. If f is a minimum of H , then, for each k , we must have $\frac{\partial H}{\partial c_k} = 0$ or, solving for c_k ,

$$c_k = \lambda_k \sum_{i=1}^m a_i \phi_k(x_i)$$

where $a_i = \frac{y_i - f(x_i)}{\gamma m}$. Thus,

$$\begin{aligned} f(x) &= \sum_{k=1}^{\infty} c_k \phi_k(x) = \sum_{k=1}^{\infty} \lambda_k \sum_{i=1}^m a_i \phi_k(x_i) \phi_k(x) \\ &= \sum_{i=1}^m a_i \sum_{k=1}^{\infty} \lambda_k \phi_k(x_i) \phi_k(x) = \sum_{i=1}^m a_i K(x_i, x). \end{aligned}$$

Replacing $f(x)$ in the definition of a_i above, we obtain

$$a_i = \frac{y_i - \sum_{j=1}^m a_j K(x_j, x)}{\gamma m}.$$

Multiplying both sides by γm and writing the result in matrix form, we obtain $(\gamma m \text{Id} + K[\mathbf{x}])a = \mathbf{y}$. And this system is well-posed since $K[\mathbf{x}]$ is positive and the addition of a positive matrix and the identity is strictly positive. \square

Classical kernel machines are equivalent to shallow networks

Kernel machines...

$$f(\mathbf{x}) = \sum_i^l c_i K(\mathbf{x}, \mathbf{x}_i) + b$$

can be “written” as shallow networks: the value of K corresponds to the “activity” of the “unit” for the input and the c_i correspond to “weights”

