

Chap 11 - Regression

1 Motivation

- (1) So far we have looked at classification problems where the results of hypotheses do not have any interesting structures — their domains are just sets, typically finite sets such as {0, 1}. The lack of an assumed structure is reflected in the definition of risk, which treats all outputs (except a single target) the same and penalize them equally. That is, the risk uses $\sum_{y' \neq y} \|h(x) - y'\|_1$, and all y' with $y' \neq y$ are equally bad as the output of $h(x)$.
- (2) But there are problems where the output domain has a useful structure, in particular, a metric, and this structure lets us say when one incorrect answer is actually better than another incorrect answer. In this chapter, we will study most well-known such problems, called regression problem.
- (3) The regression problem asks us to learn a good function from \mathcal{X} to \mathcal{Y} for $\mathcal{Y} \subseteq \mathbb{R}$ from data. The key aspect of the problem is that \mathcal{Y} has a notion of distance (i.e. $|y - y'|$), and we consider it when formalizing what we mean by good func. for a regression problem. even when $h(x) \neq y$ and $h'(x) \neq y$, we say that h is better than h' for the example (x, y) if $|h(x) - y| < |h'(x) - y|$. The objective of learning is then designed so as to reflect this changed and less rigid attitude towards incorrect outputs.

More generally - we may say that a learning problem is a regression task if the output space \mathcal{Y} is a metric space and the notion of error used for learning takes this metric into account.

(+) In this chapter, we will study the learning theory for regression and well-known algorithms for regression.

2. Problem of Regression

- (1) Setup: \mathcal{X} ... space for the inputs, $\mathcal{X} \subseteq \mathbb{R}^N$
 \mathcal{Y} ... space for the outputs, $\mathcal{Y} \subseteq \mathbb{R}$
- $D \in \Pr(\mathcal{X} \times \mathcal{Y})$ distribution that expresses what to learn
- $\mathcal{H} \subseteq \{h: \mathcal{X} \rightarrow \mathcal{Y} \mid h \text{ is measurable}\}$ hypothesis set.

(2) Learning problem

- ① Input: $S \sim D^m$, $S = ((x_1, y_1), \dots, (x_m, y_m))$
drawn independently from D .

- ② Loss function: $L: \mathcal{Y} \times \mathcal{Y} \rightarrow [0, \infty)$
- e.g., $L(y, y') = |y - y'|^2$ or $L(y, y') = |y - y'|^p$
for $p \geq 1$.

- ③ Ultimate goal: Using S , find $h \in \mathcal{H}$ that minimizes or almost minimizes. $R(h) = \mathbb{E}_{(x,y) \sim D} [L(h(x), y)]$

As we explained before, using this R instead of the one for classification is what makes regression very different from classification. By using L , R here can differentiate h, h' even when they have the same generalization error. In

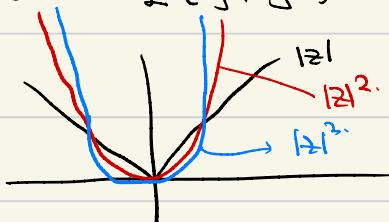
the sense of classification: $\mathbb{E} \left[\frac{1}{m} \sum_{i=1}^m \mathbb{I}_{\{h(x_i) \neq y_i\}} \right] = \mathbb{E} \left[\frac{1}{m} \sum_{i=1}^m \mathbb{I}_{\{h'(x_i) \neq y_i\}} \right]$

- ④ Most regression algorithms use the following empirical counterpart of R since the latter cannot be computed:

$$\hat{R}_S(h) = \frac{1}{m} \sum_{i=1}^m L(h(x_i), y_i)$$

- ⑤ Assume that $L(y, y') = |y - y'|^p$ for $p \geq 1$.

Note



So, using big p amounts to the decision that we penalize large deviation from the target y more, and small deviation from y less.

3. Generalisation Bounds for Bounded Regression Problems.

- (1) Before studying concrete regression algorithms, we look at generalisation bounds for bounded regression problems. Results we get here will be used to analyse the algorithms that we study later.

- (2) Let's start with the case of a finite hypothesis set, one of the simplest cases. In the context of learning theory.

Thm 11.1 $\exists M > 0$ s.t. $L(y, y') \leq M$ for all $y, y' \in \mathcal{Y}$.
Assume that \mathcal{H} is finite.

$$D \in \mathbb{P}(\mathcal{X} \times \mathcal{Y}) \quad S > 0 \quad m \in \mathbb{N}.$$

\Rightarrow

$$\underset{\text{SVD}}{\mathbb{P}} \left[\forall h \in \mathcal{H}, R(h) \leq \widehat{R}_S(h) + M \sqrt{\frac{\log |\mathcal{H}| + \log \frac{1}{\delta}}{2m}} \right] \geq 1 - \delta.$$

The thm is very similar to Thm 2.13. ($\frac{2}{\delta}$ there should be $\frac{1}{\delta}$).
except for the presence of M . Both thms have essentially the same proof.
(due to the fact that the loss at (x, y) can be at most M , instead of 1).

Prf. Let $\delta \stackrel{\text{def}}{=} M \sqrt{\frac{1}{2m} (\log |\mathcal{H}| + \log \frac{1}{\delta})}$

$$\underset{\text{SVD}}{\mathbb{P}} \left[\exists h \in \mathcal{H}, R(h) - \widehat{R}_S(h) > \delta \right].$$

$$\leq \sum_{h \in \mathcal{H}} \underset{\text{SVD}}{\mathbb{P}} \left[\sum_{i=1}^m L(h(x_i), y_i) - \left(m \mathbb{E}[L(h(x), y)] \right) < -m\delta \right].$$

Union-bound

$$\leq \sum_{h \in \mathcal{H}} \exp \left(- \frac{2(m^2 \delta^2)}{\sum_{i=1}^m M^2} \right) = \sum_{h \in \mathcal{H}} \exp \left(- \frac{2m}{M^2} \times \left(\frac{M^2}{2m} \log \left(\frac{|\mathcal{H}|}{\delta} \right) \right) \right) = \sum_{h \in \mathcal{H}} \frac{\delta}{|\mathcal{H}|} = \delta.$$

Hoeffding

D.

As in the case of classification, this bound suggests the use of the simple hypo. set. (Occam's razor principle).

(?) Usually, though, \mathcal{H} is not finite. In those cases, we will have to try what we will present next, which is derived from the Rademacher complexity bound of Thm 3.3, or other more powerful tools from learning theory.

Prop 11.2.

$L(-, y) : \mathbb{Y} \rightarrow [0, \infty)$ is μ -Lipschitz for $\mu \circ$.

$$S = ((x_1, y_1), \dots, (x_m, y_m))$$

$$\mathcal{G} = \{ (x, y) \mapsto L(h(x), y) \mid h \in \mathcal{H} \}.$$

\Rightarrow

$$\widehat{R}_S(g) \leq \mu \widehat{R}_{S^{\text{fl}}}(f)$$

$$= (x_1, \dots, x_m).$$

Proof.

$$\widehat{R}_S(g) = \frac{1}{m} \mathbb{E} \left[\sup_{\sigma \in \text{Unif}\{-1, 1\}^m} \sum_{i=1}^m g_i L(h(x_i), y_i) \right]$$

$$\leq \frac{\mu}{m} \mathbb{E} \left[\sup_{h \in \mathcal{H}} \sum_{i=1}^m g_i h(x_i) \right] = \mu \widehat{R}_{S^{\text{fl}}}(f).$$

$$\uparrow \quad \exists \tilde{a}(-) = L(-, y_{\tilde{a}}), \quad b = \mu$$

Q.

Lemma 7. [Talagrand's lemma].

$\Psi_1, \dots, \Psi_m : \mathbb{R} \rightarrow \mathbb{R}$, all ℓ -Lipschitz.

$\mathcal{H} \subseteq [\mathbb{X} \rightarrow \mathbb{R}]$,
measurable

\Rightarrow

$$\begin{aligned} & \frac{1}{m} \mathbb{E} \left[\sup_{\sigma \in \text{Unif}\{-1, 1\}^m} \left[\sup_{h \in \mathcal{H}} \sum_{i=1}^m g_i (\Psi_i(h(x_i))) \right] \right] \\ & \leq \frac{\ell}{m} \mathbb{E} \left[\sup_{h \in \mathcal{H}} \sum_{i=1}^m g_i h(x_i) \right]. \end{aligned}$$

Thm 11.3

$L: \mathcal{Y} \times \mathcal{Y} \rightarrow [0, M]$. $L(\cdot, y)$ is μ -Lipschitz.

$D \in \mathbb{P}(\mathcal{X} \times \mathcal{Y})$. $\delta > 0$. $m \in \mathbb{N}$.

\Rightarrow

$$\mathbb{P}_{\text{sup } D^m} [\forall h \in \mathcal{H}. R(h) \leq \hat{R}_S(h) + 2\mu R_m(\mathcal{H}) + M \sqrt{\frac{\log^{1/\delta}}{2m}}] \geq 1 - \delta.$$

and

$$\mathbb{P}_{\text{sup } D^m} [\forall h \in \mathcal{H}. R(h) \leq \hat{R}_S(h) + 2\mu \hat{R}_S(\mathcal{H}) + 3M \sqrt{\frac{\log^{2/\delta}}{2m}}] \geq 1 - \delta.$$

Proof. We will show only the first. Let

$$G' = \{g(x, y) \mapsto \frac{1}{m} \sum_{i=1}^m g(x_i, y_i) \mid g \in \mathcal{G}\}.$$

Then, $G' \subseteq \mathcal{F}(\mathcal{X} \times \mathcal{Y}) \rightarrow [0, 1]$, so that we can apply.

Thm 3.3. If we do so, we get:

$$\mathbb{P}_{\text{sup } D^m} [\forall g \in G'. \mathbb{E}[g(x, y)] \leq \frac{1}{m} \sum_{i=1}^m g(x_i, y_i) + 2R_m(G') + M \sqrt{\frac{\log^{1/\delta}}{2m}}] \geq 1 - \delta.$$

Since $G' = \frac{1}{m} G$ (i.e. $\{ \frac{1}{m} g \mid g \in G \}$), the above implies:

$$\mathbb{P}_{\text{sup } D^m} [\forall g \in G. \mathbb{E}[g(x, y)] \leq \frac{1}{m} \sum_{i=1}^m g(x_i, y_i) + 2M R_m(G') + M \sqrt{\frac{\log^{1/\delta}}{2m}}] \geq 1 - \delta.$$

But. $R_m(G') \leq \frac{1}{M} R_m(G)$ by Talagrand's lemma
 $\leq \frac{\mu}{M} R_m(\mathcal{H})$ by Prop 11.2.

So,

$$\mathbb{P}_{\text{sup } D^m} [\forall h \in \mathcal{H}. R(h) \leq R_S(h) + 2\mu R_m(\mathcal{H}) + M \sqrt{\frac{\log^{1/\delta}}{2m}}] \geq 1 - \delta.$$

* Without the bound M , this proof doesn't work.

* Example: $L(y, y') = |y - y'|^p$ and $|y - y'| \leq M$. for all $y, y' \in \mathcal{Y}$.

Then, $M = M_0^p$ and. $\mu = pM_0^{p-1}$ ($\because \frac{d}{dy} |y - y'|^p \leq pM_0^{p-1}$)

4. Regression algorithm 1 - Linear Regression

Linear regression is

- (1) the most basic algorithm for regression. It finds a linear fn that minimizes the empirical mean squared error.

In a sense, the algorithm is similar to ERM. for classification

- (2) Hypothesis set used: $\{f = f: \mathcal{X} \rightarrow \langle \omega, \Phi(x) \rangle + b \mid \omega \in \mathbb{R}^N, b \in \mathbb{R}\}$,
where Φ is a map from \mathcal{X} to \mathbb{R}^N and intuitively,
it computes N features for a given input x .

(3) Loss fn: $L(y, y') = (y - y')^2$

(3) Linear regression solves the following optimization problem.

$$\min_{h \in \mathcal{H}} \hat{R}_S(h)$$

where $S = ((x_1, y_1), \dots, (x_m, y_m))$ is a sample.

(or a training set).

Note:

$$= \min_{\omega, b} \frac{1}{m} \sum_{i=1}^m (\langle \omega, \Phi(x_i) \rangle + b - y_i)^2$$

$$= \min_{\omega} F(\omega)$$

where

$$X = \begin{bmatrix} \Phi(x_1) & \Phi(x_2) & \dots & \Phi(x_m) \\ 1 & 1 & \dots & 1 \end{bmatrix},$$

$$\omega = \begin{bmatrix} \omega_1 \\ \vdots \\ \omega_N \\ b \end{bmatrix}, \quad Y = \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix},$$

$$F(\omega) = \frac{1}{m} \|X^\top \omega - Y\|^2.$$

(4) The optimal ω can be obtained analytically. To see this, note that $F(\omega)$ is differentiable and convex. So, F admits a global minimum at ω iff. $\nabla F(\omega) = 0$.

But. $\nabla F(\omega) = \frac{2}{m} \times (x^T \omega - Y) = 0$. is equivalent to.
 $x x^T \omega = x Y$.

Thus, the optimal ω can be characterized as follows:

- if $(x x^T)$ is invertible, $\omega = (x x^T)^{-1} x Y$ for some ω_0 .
 - if $(x x^T)$ is not invertible,
- $$\omega = (x x^T)^+ x Y + (I - (x x^T)^+ (x x^T)) \omega_0.$$

① Here $(x x^T)^+$ is the pseudo inverse of $x x^T$. It is

$$U \Sigma^+ U^T$$

when $U \Sigma U^T$ is the eigen-decomposition of $x x^T$

and Σ^+ is the diagonal matrix with

$$(\Sigma^+)_{ii} = \begin{cases} 0 & \text{if } \Sigma_{ii} = 0 \\ \frac{1}{\Sigma_{ii}} & \text{if } \Sigma_{ii} \neq 0 \end{cases}$$

More generally, A^+ is the unique matrix s.t.

$$(i) A^+ A A^+ = A^+ \quad (ii) A A^+ A = A \quad (iii) (A A^+)^T = A^+ A$$

non-invertible.

② We can check $(x x^T) \omega = x Y$ in the case:

$$\begin{aligned} (x x^T) \omega &= (x x^T) (x x^T)^+ x Y + (x x^T) (I - (x x^T)^+ (x x^T)) \omega_0. \\ &= x x^T (x^T)^+ x^T x Y + (x x^T - (x x^T)^+ (x x^T)) \omega_0. \\ &= x (x^T (x^T)^+ (x^T x)) Y + ((x x^T) - (x x^T)^+ (x x^T)) \omega_0. \\ &= x (x^T x)^T (x^T x) Y = x (x^T x)^T x Y = x x^T x Y = x Y. \end{aligned}$$

③ In the non-muertible case, we often use:

$$\omega = (XX^T)^+ X Y$$

that has the minimum norm.

(5) ① $O(nN^2 + N^3)$ time complexity.

$$XX^T \xrightarrow{\text{inversion}}$$

② Not so good at generalisation, partly due to the lack of any effort for reducing model complexity during training.

5. Regression Algorithm 2 — Kernel Ridge Regression

(i) How should we try to reduce model complexity? The next slide gives a hint, and following the hint leads to kernel ridge regression (KRR). Assume that $L(y, y') = (y - y')^2$.

Thm 11.11

$k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$... PDS kernel

$\mathcal{H}, \mathbb{E} \dots$ RKHS, feature map of k .

$$\text{eff} = \{ \omega \mapsto \langle \omega, \mathbb{E}(x) \rangle \mid \|\omega\|_{\mathcal{H}} \leq \Lambda \}$$

$$\exists r, M \text{ s.t. } K(x, x) \leq r^2 \text{ for all } x \text{ and.}$$

$$|h(x) - y| \leq M \text{ for all } h \in \text{eff}, y \in \mathcal{Y}.$$

$D \in \mathcal{P}(\mathcal{X} \times \mathcal{Y})$, $\delta > 0$, $m \in \mathbb{N}$.

\Rightarrow

$$\underset{S \sim D^m}{\mathbb{P}} [\forall h \in \text{eff. } R(h) \leq \hat{R}_S(h) + 4M \sqrt{\frac{r^2 \Lambda^2}{m} + M^2} \sqrt{\frac{\log \frac{1}{\delta}}{2m}}] \geq 1 - \delta.$$

and

$$\underset{S \sim D^m}{\mathbb{P}} [\forall h \in \text{eff. } R(h) \leq \hat{R}_S(h) + 4M \sqrt{\frac{\text{Tr}[K]}{m}} + 3M^2 \sqrt{\frac{\log \frac{2}{\delta}}{2m}}] \geq 1 - \delta.$$

where $K = [k(x_i, x_j)]_{ij}$.

Prf. Immediate consequence of Thm 3 and Thm 12.

$$\widehat{R}_S(\delta) \leq \frac{\sqrt{\lambda \text{Tr}[K]}}{m} \leq \sqrt{\frac{\lambda^2 r^2}{m}}$$

(2) One lesson of the thm is that one way to reduce model complexity is to use small λ . We can turn this lesson to an algorithm, which becomes KRR.

KRR solves the following optimization pb:

$$\min_{\omega} F(\omega)$$

where:

$$F(\omega) = \lambda \|\omega\|^2 + \sum_{i=1}^m (\langle \omega, \Phi(x_i) \rangle - y_i)^2 \\ = \lambda \|\omega\|^2 + \|X^T \omega - Y\|^2,$$

$$X = \begin{bmatrix} \Phi(x_1) & \dots & \Phi(x_m) \end{bmatrix}, Y = \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix}.$$

(3) The solution of this problem can be computed analytically.

Since F is convex and differentiable, F attains a global minimum at ω iff $\nabla F(\omega) = 0$. But

$$\nabla F(\omega) = \lambda \omega + X(X^T \omega - Y).$$

so that $\nabla F(\omega) = 0$ is equivalent to

$$(X^T X + \lambda I) \omega = X^T Y.$$

the sums of

$X^T X + \lambda I$ is invertible because its eigenvalues are those of $X^T X$, which are ≥ 0 , and $\lambda > 0$ and so they are all positive. (The determinant is the product of eigenvalues). Thus, the optimal ω satisfies

$$\omega = (X^T X + \lambda I)^{-1} X^T Y.$$

(4) To obtain a kernelized variant of the above algorithm, we consider a related optimization problem.

$$\min_{\omega} \sum_{i=1}^m (\langle \omega, \Phi(x_i) \rangle - y_i)^2$$

s.t. $\|\omega\|^2 \leq \lambda^2$

=

$$\min_{\tilde{z}, \omega} \sum_{i=1}^m \tilde{z}_i^2$$

the condition $d \geq 0$ because d is introduced for the equality constraint

$$\text{s.t. } \|\omega\|^2 \leq \lambda^2 \wedge \tilde{z}_i = y_i - \langle \omega, \Phi(x_i) \rangle \text{ for } i \in \mathbb{R}^m$$

where

$$\min_{\tilde{z}, \omega} \max_{d' \geq 0, \lambda \geq 0} \mathcal{L}(\tilde{z}, \omega, d', \lambda).$$

$$\mathcal{L}(\tilde{z}, \omega, d', \lambda) = \sum_{i=1}^m \tilde{z}_i^2 + \sum_{i=1}^m d'_i (y_i - \langle \omega, \Phi(x_i) \rangle - \tilde{z}_i) + \lambda (\|\omega\|^2 - \lambda^2)$$

=

$$\max_{d' \geq 0, \lambda \geq 0} \min_{\tilde{z}, \omega} \mathcal{L}(\tilde{z}, \omega, d', \lambda)$$

KKT condition:

$$\nabla_{\omega} \mathcal{L} = 0 \Rightarrow \omega = \frac{1}{2\lambda} \sum_{j=1}^m d_j' \Phi(x_j)$$

$$\nabla_{\tilde{z}} \mathcal{L} = 0 \Rightarrow \tilde{z}_i = d_i' / 2.$$

$$\forall i \in \mathbb{R}^m, d_i' (y_i - \langle \omega, \Phi(x_i) \rangle - \tilde{z}_i) = 0.$$

$$\lambda (\|\omega\|^2 - \lambda^2) = 0.$$

$$\begin{aligned} &= \max_{d' \geq 0, \lambda \geq 0} \sum_{i=1}^m \frac{d_i'^2}{4} + \sum_{i=1}^m d_i' (y_i - \frac{1}{2\lambda} \sum_{j=1}^m d_j' k(x_j, x_i) - \frac{d_i'}{2}) \\ &\quad + \frac{1}{4\lambda} \sum_{i=1}^m \sum_{j=1}^m d_i' d_j' k(x_i, x_j) - \lambda \lambda^2 \end{aligned}$$

$$= \max_{d^*, \lambda \geq 0} -\frac{1}{4} \sum_{i=1}^m d_i'^2 + \sum_{i=1}^m d_i y_i - \frac{1}{4} \lambda \sum_{i,j=1}^m d_i d_j k(x_i, x_j) - \lambda \lambda^2$$

Let $d_i' = 2\sqrt{\lambda} d_i$

$$\begin{aligned} &= \max_{d, \lambda \geq 0} -\lambda \sum_{i=1}^m d_i^2 + 2\sqrt{\lambda} \sum_{i=1}^m d_i y_i - \sum_{i,j=1}^m d_i d_j k(x_i, x_j) - \lambda \lambda^2 \\ &= \max_{d, \lambda \geq 0} -\lambda d^T d + 2\sqrt{\lambda} d^T Y - d^T K d - \lambda \lambda^2 \end{aligned}$$

(and ignore $\rightarrow \lambda^2$)

Now if we treat λ as a given constant, this opt. is equivalent.
Thus, the kernelized dual version of KRR solver:

$$\operatorname{argmax}_{d \in \mathbb{R}^m} -\lambda d^T d + 2\sqrt{\lambda} d^T Y - d^T K d.$$

$= G(d).$

G is concave and differentiable. Thus, it has a global optimum.

at d . iff $\nabla_d G(d) = 0$.

$$\text{But } \nabla_d G(d) = -2\lambda d + 2\sqrt{\lambda} Y - 2Kd.$$

Thus, $\nabla_d G(d) = 0$ is equivalent to.

$$d = \sqrt{\lambda} (\lambda I + K)^{-1} Y$$

In this case,

$$\begin{aligned} \langle w, \Phi(x) \rangle &= \frac{1}{2\lambda} \sum_{i=1}^m 2\sqrt{\lambda} d_i \langle \Phi(x_i), \Phi(x) \rangle = \frac{1}{\sqrt{\lambda}} \sum_{i=1}^m d_i k(x_i, x) \\ &= \frac{1}{\sqrt{\lambda}} K_{x, x_{1:m}} d = \frac{1}{\sqrt{\lambda}} K_{x, x_{1:m}} \sqrt{\lambda} (\lambda I + K)^{-1} Y \\ &= K_{x, x_{1:m}} (\lambda I + K)^{-1} Y. \end{aligned}$$

almost.

(5) KRR with bias term b . can be covered by what we discussed so far if we use $\Phi^l(x) = \begin{bmatrix} \Phi(x) \\ 1 \end{bmatrix} \in \mathbb{R}^{N+1}$.
But in this case, we optimise $\lambda \| [w_b] \|^2 + \dots$ instead of $\lambda \| w \|^2 + \dots$.

(b) Computational Complexity

m examples.

N features.

K = cost of computing $k(x, x')$.

	Learning	Prediction
primal	$O(mN^2 + N^3)$	$O(N)$
dual	$O(km^2 + m^3)$	$O(mk)$

So, depending on whether $N > m$ or $N < m$, we may prefer primal version or dual version.

(c) One drawback of KRR is that its solution is not sparse. The next two algorithms are designed to address this issue.

b. Regression Algorithm 3 - Support Vector Regression

(SVR)

(1) The support vector regression is a regression algorithm that often finds a sparse solution, i.e., a fn that uses a small number of features or dimensions in the input. The key idea of the SVR is not to penalize a prediction if it is close enough (ϵ -close) to the ground truth. It lets the SVR achieve sparsity and be analysed using concepts from the SVM, such as support vector.

(2) The SVR solves the following optimization problem:

① Input : $S = ((x_1, y_1), \dots, (x_m, y_m))$

Hyperparameters : $C > 0$, $\epsilon > 0$.

$$\textcircled{2} \quad \min_{\omega, b} \frac{1}{2} \|\omega\|^2 + C \sum_{i=1}^m |y_i - (\langle \omega, \Phi(x_i) \rangle + b)|_\varepsilon$$

where $|y|_\varepsilon = \max(0, |y|_\delta)$

feature map

=

$$\min_{\omega, b, \beta, \beta'} \frac{1}{2} \|\omega\|^2 + C \sum_{i=1}^m (\beta_i + \beta'_i)$$

subject to. $\beta'_i \geq 0 \quad \wedge \quad \beta'_i \geq y_i - (\langle \omega, \Phi(x_i) \rangle + b) - \varepsilon$
 $\beta_i \geq 0 \quad \wedge \quad \beta_i \geq (\langle \omega, \Phi(x_i) \rangle + b) - y_i - \varepsilon$

for all $i \in [m]$.

Convex quadratic optimisation pb. with linear constraints.

Can be solved by an off-the-shelf optimiser.

(3) The optimisation from above can be dualised as shown below, and can be converted into a kernelised algorithm.

$$\textcircled{1} \quad \mathcal{L}(\omega, b, \beta, \beta', \alpha, \alpha', \beta - \beta')$$

$$= \frac{1}{2} \|\omega\|^2 + C \sum_{i=1}^m (\beta_i + \beta'_i) + \sum_{i=1}^m (\beta_i (-\beta'_i) + \beta'_i (-\beta_i))$$

$$+ \sum_{i=1}^m (\alpha_i ((\langle \omega, \Phi(x_i) \rangle + b) - y_i - \varepsilon) + \alpha'_i (y_i - (\langle \omega, \Phi(x_i) \rangle + b) - \varepsilon))$$

$$\min_{\omega, b, \beta, \beta'} \dots = \inf_{\omega, b, \beta, \beta'} \sup_{\alpha, \alpha' \in \mathbb{R}_0, \infty} \mathcal{L}(\omega, b, \beta, \beta', \alpha, \alpha', \beta - \beta')$$

$$= \sup_{\alpha, \alpha' \in \mathbb{R}_0, \infty} \inf_{\omega, b, \beta, \beta'} \mathcal{L}(\omega, b, \beta, \beta', \alpha, \alpha', \beta - \beta').$$

KKT cond. $\nabla_\omega \mathcal{L} = 0 \Rightarrow \omega = -\sum_{i=1}^m (\alpha_i - \alpha'_i) \Phi(x_i)$.

$$\nabla_b \mathcal{L} = 0 \Rightarrow \sum_{i=1}^m (\alpha_i - \alpha'_i) = 0. \quad \nabla_{\beta, \beta'} \mathcal{L} = 0 \Rightarrow C = \alpha_i + \alpha'_i.$$

$$\nabla_{\beta, \beta'} \mathcal{L} = 0 \Rightarrow C = \alpha'_i + \beta'_i.$$

Let $k(x, x') = \langle \Phi(x), \Phi(x') \rangle$ and $K_{ij} = k(x_i, x_j)$
 $\sup \inf \{ \dots \text{ from above.} \}$ \hookrightarrow kernel or Gram matrix

$$= \sup_{d, d'} -\frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m (d_i - d'_i)(d_j - d'_j) \langle \Phi(x_i), \Phi(x_j) \rangle - \sum_{i=1}^m (d_i - d'_i) y_i - \varepsilon \sum_{i=1}^m (d_i + d'_i)$$

s.t. $0 \leq d_i \leq C \wedge 0 \leq d'_i \leq C \wedge \sum_{i=1}^m (d_i - d'_i) = 0$

$$= \sup_{d, d'} -\varepsilon \langle d + d'^T, 1 \rangle - \langle d - d', y \rangle - \frac{1}{2} \langle d - d', K(d - d') \rangle$$

s.t. $0 \leq d_i \leq C \wedge 0 \leq d'_i \leq C \wedge \langle d - d', 1 \rangle = 0 \text{ for all } i \in [m]$

② Hypothesis $h^*(x) = \langle \omega^*, \Phi(x) \rangle + b$.

$$= \left\langle \sum_{i=1}^m (d_i^* - d_i'^*) \Phi(x_i), \Phi(x) \right\rangle + b^*$$

$$= \sum_{i=1}^m (d_i^* - d_i'^*) k(x_i, x) + b^*$$

For b^* , we pick \bar{i}_0 with $0 < d_{\bar{i}_0}^* < C$ or

$$0 < d_{\bar{i}_0}^{1*} < C.$$

In the former case, $\bar{x}_{\bar{i}_0} = 0$ and $\langle \omega^*, \Phi(\bar{x}_{\bar{i}_0}) \rangle + b^* - y_{\bar{i}_0} - \varepsilon = 0$.

by KKT. So, $b^* = y_{\bar{i}_0} + \varepsilon + \sum_{i=1}^m (d_i^* - d_i'^*) k(x_i, \bar{x}_{\bar{i}_0})$.

We can derive b^* similarly in the latter case.

③ Note that only the \bar{i} 's with $d_i^* - d_i'^* \neq 0$ matter for h^* . So,

x_i with such \bar{i} plays the role of support vectors in the SVM algo.

If (x_i, y_i) is ε -close to h^* , then by KKT, $d_i^* = d_i'^* = 0$.

$$(|h^*(x_i) - y_i| = |\langle \omega^*, \Phi(x_i) \rangle + b^* - y_i| < \varepsilon)$$

Thus, such input doesn't contribute to the def'n of h^* .

④ Compare this dual opt. pb. with the dual objective of KRR

$$\text{SVR: } \max_{\alpha, \alpha'} -\varepsilon (\alpha + \alpha')^T \mathbf{1} + (\alpha' - \alpha)^T \mathbf{y} - \frac{1}{2} (\alpha' - \alpha)^T K (\alpha' - \alpha) \\ \text{s.t. } 0 \leq \alpha \leq C \wedge 0 \leq \alpha' \leq C \wedge (\alpha' - \alpha)^T \mathbf{1} = 0. \\ (\alpha_i \leq C \text{ for all } i \in [m])$$

$$\text{KRR: } \max_{\alpha \in \mathbb{R}^m} \mathbf{z}^T \mathbf{y} - \alpha^T (K + \lambda I) \alpha.$$

$\alpha' - \alpha$ in SVR corresponds to α in KRR. So, the SVR objective has an extra term $-\varepsilon (\alpha + \alpha')^T \mathbf{1}$ and more constraints, but it also lacks the λI part. The missing λI part is due to the use of the L_1 norm when we compute the empirical loss. If we used the L_2 norm instead, that is, we use the following objective

$$\max_{w, b} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \|y_i - \langle w, \Phi(x_i) \rangle - b\|_2^2$$

then the dualization would give the following objective:

$$\max_{\alpha, \alpha'} -\varepsilon (\alpha + \alpha')^T \mathbf{1} + (\alpha' - \alpha)^T \mathbf{y} - \frac{1}{2} (\alpha' - \alpha)^T (K + \frac{1}{C} I) (\alpha' - \alpha) \\ \text{s.t. } 0 \leq \alpha \leq C \wedge 0 \leq \alpha' \leq C \wedge (\alpha' - \alpha)^T \mathbf{1} = 0.$$

Note that if $\varepsilon = 0$ and we ignore the $(\alpha' - \alpha)^T \mathbf{1} = 0$ constraint, which appears due to the presence of the bias term b . (in our analysis of KRR, we ignored the bias term), then the above objective becomes identical to that of KRR. (modulo the correspondence between $\alpha' - \alpha$ in SVR and α in KRR).

(4). Generalisation bounds. We ignore the bias term b here.

① SVR with the L_1 empirical loss.

Thm 11.13.

$k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$... PDS kernel. ($K_{ij} = k(x_i, x_j)$)

Gram matrix

for k and S

\mathbb{H}, Ξ RKHS, feature map of k .

$$H = \{x \mapsto \langle w, \Xi(x) \rangle_{\mathbb{H}} \mid \|w\|_{\mathbb{H}} \leq 1\}.$$

$\exists r > 0, M > 0$ s.t. $|K(x, x)| \leq r^2$ for all $x \in \mathcal{X}$.

and $|y - y'| \leq M$ for all $y, y' \in \mathcal{Y}$.

$D \in \mathbb{R}(\mathcal{X} \times \mathcal{Y})$, $\varepsilon > 0$, $\delta > 0$, $m \in \mathbb{N}$.

\Rightarrow

$$\underset{\text{sup } m}{P} \left[\text{the eff. } \mathbb{E} \sum_{(x,y) \sim D} |h(x) - y|_{\varepsilon} \right] \leq \frac{1}{m} \sum_{i=1}^m |h(x_i) - y_i|_{\varepsilon} + 2 \sqrt{\frac{r^2 \Lambda^2}{m}} + M \sqrt{\frac{\log(1/\delta)}{2m}} \geq 1 - \delta$$

$$\text{and} \quad \underset{\text{sup } m}{P} \left[\text{the eff. } \mathbb{E} \sum_{(x,y) \sim D} |h(x) - y|_{\varepsilon} \right] \leq \frac{1}{m} \sum_{i=1}^m |h(x_i) - y_i|_{\varepsilon} + \frac{2\Lambda \sqrt{\text{Tr}[K]}}{m} + 3M \sqrt{\frac{\log^2(1/\delta)}{2m}} \geq 1 - \delta$$

Pro-f. Let $L: \mathcal{Y} \times \mathcal{Y} \rightarrow [0, \infty)$ be $L(y, y') \stackrel{\text{def.}}{=} |y - y'|_{\varepsilon}$.

Then, L is 1-Lip and bounded by M . Also, by Thm 6.12,

$$\widehat{R}_S(\text{eff.}) \leq \frac{\sqrt{\text{Tr}[K]}}{m} \leq \sqrt{\frac{\Lambda^2 r^2}{m}}. \quad \text{So, by Thm 11.3,}$$

we have the bounds in the theorem. \square .

② SVR with the L_2 empirical loss. Essentially by the same reasoning, we can get the following generalization bounds.

Thm 11.14. Same setup as the one for Thm 11.13. Then,

$$\underset{\text{sup } m}{P} \left[\text{the eff. } \mathbb{E} [|h(x) - y|_{\varepsilon}^2] \leq \frac{1}{m} \sum_{i=1}^m |h(x_i) - y_i|_{\varepsilon}^2 + 4M \sqrt{\frac{\Lambda^2 r^2}{m}} + M^2 \sqrt{\frac{\log(1/\delta)}{2m}} \right] \geq 1 - \delta \text{ and.}$$

$$\underset{\text{sup } m}{P} \left[\text{the eff. } \mathbb{E} [|h(x) - y|_{\varepsilon}^2] \leq \frac{1}{m} \sum_{i=1}^m |h(x_i) - y_i|_{\varepsilon}^2 + \frac{4M \sqrt{\text{Tr}[K]}}{m} + 3M^2 \sqrt{\frac{\log^2(1/\delta)}{2m}} \right] \geq 1 - \delta.$$

Prоф. The same proof as the one for Thm 1.3 except that
 the loss fn $L(y, y') = \max(0, |y - y'| - \varepsilon)^2$ is bounded by M^2
 and is $2M$ -Lipschitz. R.

(5) Any problem with the SVR? Yes. First, we have to find appropriate values for the hyperparameters C and ε . Second, the algorithm and the learnt classifier are slow when m is large. The next algorithm, Lasso, addresses these issues at least partially, which computing sparse solutions.

7. Regression Algorithm t - Lasso.

(1) The key idea of Lasso is to use the L_1 regularization on weights, instead of L_2 . Lasso solves the following opti. pb:

$$\min_{\omega, b} \lambda \|\omega\|_1 + \sum_{i=1}^m (\langle \omega, \Phi(x_i) \rangle + b - y_i)^2$$

$\lambda \sum_{k=1}^N |\omega_k|$

Solutions of Lasso tend to be sparse. That is, $\{\omega_i \in [N]\}_{i=1}^N$ tends to be small. The reason is that the optimal values are often obtained at corners. See Fig 11.6 in the textbook.

(2) How to solve the above optimization problem? One issue is that $\|\omega\|_1$ is the sum of absolute values, and the absolute value fn is not as nice as polynomials. The next transformation provides a way to get rid of the absolute values.

$$\min_{\omega, b} \lambda \|\omega\|_1 + \sum_{i=1}^m (\langle \omega, \Phi(x_i) \rangle + b - y_i)^2$$

$$= \min_{\substack{\omega^+ \\ \omega^- \\ b}} \lambda \sum_{k=1}^n (\omega_k^+ + \omega_k^-) + \sum_{i=1}^m (\langle \omega^+ - \omega^-, \Phi(x_i) \rangle + b - y_i)^2$$

$$\omega^+ \geq 0, \omega^- \geq 0, b.$$

The new optimization problem is an instance of so called quadratic programming, and it can be solved by an off-the-shelf solver.

(3) We can use a generalization bound for Lasso again using Thm 11.3. But this time we want to bound the Rademacher complexity part in the conclusion of Thm 11.3 differently. We want the L_1 norm to feature there, instead of the L_2 norm.

The next time is crucial to achieve what we want.

By the way, in this theoretical analysis of Lasso, we ignore the bias term b .

Thm 11.15:

$$S = (x_1, \dots, x_m) \in \mathbb{R}^m$$

$$\exists r_\infty > 0 \text{ s.t. } \forall i \in [m]. \quad \|\Phi(x_i)\|_\infty < r_\infty$$

$$\max_{j \in [m]} |\Phi(x_j)_j|.$$

$$f = \{x \mapsto \langle \omega, \Phi(x) \rangle \mid \|\omega\|_1 \leq \Lambda, \}$$

\Rightarrow

$$\widehat{R}_S(f) \leq \sqrt{\frac{2r_\infty^2 \Lambda^2 \log(2N)}{m}}$$

$$\begin{aligned} \text{Proof. } \widehat{R}_S(\mathcal{H}) &= \frac{1}{m} \mathbb{E}_{\vec{\omega} \sim \text{Unif}\mathbb{S}^{m-1}, \mathcal{B}^m} \left[\sup_{h \in \mathcal{H}} \sum_{i=1}^m g_i h(x_i) \right] \\ &= \frac{1}{m} \mathbb{E}_{\vec{\omega}} \left[\sup_{\|\omega\|_1 \leq \Lambda_1} \sum_{i=1}^m g_i \langle \omega, \Phi(x_i) \rangle \right] \end{aligned}$$

$$\begin{aligned} &= \frac{1}{m} \mathbb{E}_{\vec{\omega}} \left[\sup_{\|\omega\|_1 \leq \Lambda_1} \left\langle \omega, \sum_{i=1}^m g_i \Phi(x_i) \right\rangle \right] \\ &\leq \frac{1}{m} \mathbb{E}_{\vec{\omega}} \left[\sup_{\|\omega\|_1 \leq \Lambda_1} \|\omega\|_1 \left\| \sum_{i=1}^m g_i \Phi(x_i) \right\|_\infty \right]. \end{aligned}$$

$$= \frac{\Lambda_1}{m} \mathbb{E}_{\vec{\omega}} \left[\max_{j \in [N]} \left| \sum_{i=1}^m g_i \Phi(x_i)_j \right| \right]$$

$$= \frac{\Lambda_1}{m} \mathbb{E}_{\vec{\omega}} \left[\max_{j \in [N]} \max_{s \in \{-1, +1\}} \sum_{i=1}^m g_i (s \Phi(x_i)_j) \right]$$

Let

$$A = \left\{ (s \Phi(x_1)_j, s \Phi(x_2)_j, \dots, s \Phi(x_m)_j) \mid s \in \{-1, +1\}, j \in [N] \right\}$$

$$= \frac{\Lambda_1}{m} \mathbb{E}_{\vec{\omega}} \left[\sup_{z \in A} \sum_{i=1}^m g_i z_i \right]$$

Since for any $z \in A$, $\|z\|_2 \leq \sqrt{m r_\infty^2}$ and $|A| \leq 2N$,

by Massart's lemma (Thm 3.11),

$$\leq \frac{\Lambda_1}{m} \cdot \sqrt{m r_\infty^2} \sqrt{2 \log 2N} = \Lambda_1 r_\infty \sqrt{\frac{2 \log 2N}{m}} \quad \square$$

Thm 3.11. $A \subseteq \mathbb{R}^m$ finite set. $r = \max_{x \in A} \|x\|_2$.

$$\Rightarrow \mathbb{E}_{\vec{\omega} \sim \text{Unif}\mathbb{S}^{m-1}, \mathcal{B}^m} \left[\frac{1}{m} \sup_{x \in A} \sum_{i=1}^m g_i x_i \right] \leq r/m \times \sqrt{2 \log |A|} //$$

Thm 11.16

$$\mathcal{H} = \{x \mapsto \langle \omega, \Phi(x) \rangle \mid \|\omega\|_1 \leq \Lambda_1\}.$$

$$\exists r_\infty > 0 \text{ s.t. } \forall x \in \mathcal{X} \quad \|x\|_\infty \leq r_\infty.$$

$$\exists M > 0 \text{ s.t. } \forall y, y' \in \mathcal{Y} \quad |y - y'| \leq M.$$

$$D \in \mathbb{R}^{(X \times Y)} \quad m \in \mathbb{N} \quad \delta > 0.$$

\Rightarrow

$$\underset{S \sim D^m}{P} \left[\text{The eff. Rch} \right] \leq \widehat{R}_{\text{SCh}} + \frac{\tau_\infty \Lambda_1 M \sqrt{\frac{2 \log 2N}{m}}}{+ M^2 \sqrt{\frac{\log 1/\delta}{2m}}} \geq 1 - \delta.$$

$\widehat{R}_{\text{SCh}} = \mathbb{E}_{(x,y) \sim D} [C_h(x) - y]^2$

Pf. By Thm 11.3 and Thm 11.15. We use the fact that

the loss fn $L(y, y') = (y - y')^2$ is bounded by M^2 and
and $2M$ -Lipschitz. \square

(4) Any drawback for Lasso? First, it doesn't permit a natural use of PDS kernels. So, kernelizing it and generalizing it to the non-linear case are hard. Second, it doesn't admit a closed-form solution. This is not an issue in practice, but it makes the theoretical analysis of Lasso difficult.