The Learning Problem and Regularization

Tomaso Poggio

9.520 Class 02

September 2015

Tomaso Poggio The Learning Problem and Regularization

・ 同 ト ・ ヨ ト ・ ヨ ト

ъ

Statistical Learning Theory

Learning is viewed as a generalization/inference problem from usually **small** sets of **high dimensional**, noisy data.

Today's class is one of the most difficult – because it is abstract. Reasons for it:

- Science of Learning
- Big picture and flavor
- Mathcamp is next
- This classroom is not large enough.

・ 同 ト ・ ヨ ト ・ ヨ

There are in principle several "learning problems". The one which is most crisply defined is supervised learning. If the conjecture about Implicit Supervised Examples were correct, then *supervised learning* – together with reinforcement learning – would be the most important building block for the whole of biological learning.

- Supervised
- Semisupervised
- Unsupervised
- Online
- Transductive
- Active
- Variable Selection
- Reinforcement
-

In addition one can consider the data to be created in a deterministic, or stochastic or even adversarial way.

くロト (過) (目) (日)

Where to Start?

Statistical and Supervised Learning

- Statistical Models are essentially to deal with noise sampling and other sources of uncertainty.
- Supervised Learning is the best understood type of learning problems and may be a building block for most of the others.

Regularization

- Regularization provides a rigorous framework to solve learning problems and to design learning algorithms.
- In the course we will present a set of ideas and tools which are at the core of several developments in supervised learning and beyond it.

We will see the close connection during the last classes between kernel machines and deep networks.

ヘロア 人間 アメヨア 人口 ア

Where to Start?

Statistical and Supervised Learning

- Statistical Models are essentially to deal with noise sampling and other sources of uncertainty.
- Supervised Learning is the best understood type of learning problems and may be a building block for most of the others.

Regularization

- Regularization provides a rigorous framework to solve learning problems and to design learning algorithms.
- In the course we will present a set of ideas and tools which are at the core of several developments in supervised learning and beyond it.

We will see the close connection during the last classes between kernel machines and deep networks.

ヘロト 人間 ト ヘヨト ヘヨト

Remarks on Foundations of Learning Theory

- This class establish our program for the first 10 classes:
 - Main goal of learning is generalization and predictivity not explanation
 - Which algorithms to guarantee ensure generalization?
 - We derive "equivalence" of generalization and stability/well-posedness
 - Since it is known that regularization techniques guarantee well-posedness we will use them to guarantee also generalization
 - Notice that they usually result in *computationally "nice"* and *well-posed* constrained optimization problems

・ロット (雪) () () () ()

Part I: Basic Concepts and Notation

- Part II: Foundational Results
- Part III: Algorithms

・ 同 ト ・ ヨ ト ・ ヨ ト

ъ

Given a training set of input-output pairs

$$S_n = (x_1, y_1), \ldots, (x_n, y_n)$$

find f_S such that

$$f_S(x) \sim y$$
.

e.g. the x's are vectors and the y's discrete labels in classification and real values in regression.

(同) くほり くほう

For the above problem to make sense we need to assume input and output to be related!

Statistical and Supervised Learning

- Each input-output pairs is a sample from a fixed but unknown distribution μ(x, y).
- Under some condition we can associate to $\mu(z)$ the probability

p(x,y) = p(y|x)p(x).

- the training set S_n is a set of identically and independently distributed samples drawn from μ(z).
- It is crucial to note that we view p(x, y) as fixed but unknown.

・ロ・ ・ 四・ ・ ヨ・ ・ ヨ・

For the above problem to make sense we need to assume input and output to be related!

Statistical and Supervised Learning

- Each input-output pairs is a sample from a fixed but unknown distribution $\mu(x, y)$.
- Under some condition we can associate to $\mu(z)$ the probability

p(x,y)=p(y|x)p(x).

- the training set S_n is a set of identically and independently distributed samples drawn from μ(z).
- It is crucial to note that we view p(x, y) as fixed but unknown.

ヘロト ヘ戸ト ヘヨト ヘヨト

Why Probabilities



the same x can generate different y (according to p(y|x)):

- the underlying process is deterministic, but there is noise in the measurement of y;
- the underlying process is not deterministic;
- the underlying process is deterministic, but only **incomplete** information is available.



even in a noise free case we have to deal with sampling

the marginal p(x) distribution might model

- errors in the location of the input points;
- discretization error for a given grid;
- presence or absence of certain input instances

★ E > ★ E



even in a noise free case we have to deal with sampling

the marginal p(x) distribution might model

- errors in the location of the input points;
- discretization error for a given grid;
- presence or absence of certain input instances

∃ → <</p>



even in a noise free case we have to deal with sampling

the marginal p(x) distribution might model

- errors in the location of the input points;
- discretization error for a given grid;
- presence or absence of certain input instances

글 🕨 🖌 글



even in a noise free case we have to deal with sampling

the marginal p(x) distribution might model

- errors in the location of the input points;
- discretization error for a given grid;
- presence or absence of certain input instances

Given a training set of input-output pairs

$$S_n = (x_1, y_1), \ldots, (x_n, y_n)$$

find f_S such that

 $f_S(x) \sim y$.

e.g. the x's are vectors and the y's discrete labels in classification and real values in regression.

個 ト く ヨ ト く ヨ ト

Predictivity or Generalization

Given the data, the goal is to learn how to make decisions/predictions about future data / data not belonging to the training set. **Generalization** is the key requirement emphasized in Learning Theory: generalization is a masure of predictivity. This emphasis makes it different from Bayesian or traditional statistics (especially explanatory statistics).

The problem is often: Avoid overfitting!!

・ 同 ト ・ ヨ ト ・ ヨ ト

In order to define generalization we need to define and measure errors.

Loss function

A loss function $V : \mathbf{R} \times Y$ determines the price V(f(x), y) we pay, predicting f(x) when in fact the true output is y.

ヘロン 人間 とくほ とくほ とう

3

In order to define generalization we need to define and measure errors.

Loss function

A loss function $V : \mathbf{R} \times Y$ determines the price V(f(x), y) we pay, predicting f(x) when in fact the true output is y.

くロト (過) (目) (日)

The most common is the square loss or L₂ loss

$$V(f(x), y) = (f(x) - y)^2$$

• Absolute value or L₁ loss:

$$V(f(x), y) = |f(x) - y|$$

• Vapnik's *e*-insensitive loss:

$$V(f(x), y) = (|f(x) - y| - \epsilon)_+$$

Loss functions for (binary) classification

• The most intuitive one: 0 – 1-loss:

$$V(f(x), y) = \theta(-yf(x))$$

(θ is the step function)

• The more tractable hinge loss:

$$V(f(x), y) = (1 - yf(x))_+$$

• And again the square loss or L₂ loss

$$V(f(x), y) = (1 - yf(x))^2$$

Loss functions



◆□> ◆□> ◆豆> ◆豆> ・豆 ・ のへで

A good function – we will speak of function or *hypothesis* – should incur in only *a few* errors. We need a way to quantify this idea.

Expected Risk

The quantity

$$I[f] = \int_{X \times Y} V(f(x), y) p(x, y) dx dy.$$

is called the expected error and measures the loss averaged over the unknown distribution.

A good function should have small expected risk.

イロト イポト イヨト イヨト

A good function – we will speak of function or *hypothesis* – should incur in only *a few* errors. We need a way to quantify this idea.

Expected Risk

The quantity

$$I[f] = \int_{X \times Y} V(f(x), y) p(x, y) dx dy.$$

is called the expected error and measures the loss averaged over the unknown distribution.

A good function should have small expected risk.

ヘロト ヘ戸ト ヘヨト ヘヨト

A learning algorithm can be seen as a map

$$S_n \rightarrow f_n$$

from the training set to the a set of candidate functions.

- p(x, y) probability distribution,
- S_n training set,
- V(f(x), y) loss function,
- $I_n[f] = \frac{1}{n} \sum_{i=1}^n V(f(x_i), y_i)$, empirical risk,
- $I[f] = \int_{X \times Y} V(f(x), y) p(x, y) dx dy$, expected risk.

Reminder

if

Convergence in probability

Let $\{X_n\}$ be a sequence of bounded random variables. Then

 $\lim_{n\to\infty} X_n = X \quad \text{in probability}$

$$\forall \epsilon > 0 \quad \lim_{n \to \infty} \mathbb{P}\{|X_n - X| \ge \epsilon\} = 0$$

Convergence in Expectation

Let $\{X_n\}$ be a sequence of bounded random variables. Then

$$\lim_{n\to\infty} X_n = X \quad \text{in expectation}$$

$$\lim_{n\to\infty}\mathbb{E}(|X_n-X|)=0$$

Convergence in the mean implies convergence in probability.

Reminder

Convergence in probability

Let $\{X_n\}$ be a sequence of bounded random variables. Then

 $\lim_{n\to\infty} X_n = X \text{ in probability}$

if

$$\forall \epsilon > 0 \quad \lim_{n \to \infty} \mathbb{P}\{|X_n - X| \ge \epsilon\} = 0$$

Convergence in Expectation

Let $\{X_n\}$ be a sequence of bounded random variables. Then

$$\lim_{n\to\infty} X_n = X \quad \text{in expectation}$$

$$\lim_{n\to\infty}\mathbb{E}(|X_n-X|)=0$$

Convergence in the mean implies convergence in probability. Tomaso Poggio

The Learning Problem and Regularization

Consistency and Universal Consistency

A requirement considered of basic importance in classical statistics is for the algorithm to get better as we get more data (in the context of machine learning consistency is less immediately critical than *generalization*)...

Consistency

We say that an algorithm is consistent if

$$\forall \epsilon > 0 \quad \lim_{n \to \infty} \mathbb{P}\{I[f_n] - I[f_*] \ge \epsilon\} = 0$$

Universal Consistency

We say that an algorithm is universally consistent if for all probability p,

$$\forall \epsilon > 0 \quad \lim_{n \to \infty} \mathbb{P}\{I[f_n] - I[f_*] \ge \epsilon\} = 0$$

ヘロト ヘ戸ト ヘヨト ヘヨト

Consistency and Universal Consistency

A requirement considered of basic importance in classical statistics is for the algorithm to get better as we get more data (in the context of machine learning consistency is less immediately critical than *generalization*)...

Consistency

We say that an algorithm is consistent if

$$\forall \epsilon > 0 \quad \lim_{n \to \infty} \mathbb{P}\{I[f_n] - I[f_*] \ge \epsilon\} = 0$$

Universal Consistency

We say that an algorithm is universally consistent if for all probability *p*,

$$\forall \epsilon > 0 \quad \lim_{n \to \infty} \mathbb{P}\{I[f_n] - I[f_*] \ge \epsilon\} = 0$$

イロト イポト イヨト イヨト

Consistency and Universal Consistency

A requirement considered of basic importance in classical statistics is for the algorithm to get better as we get more data (in the context of machine learning consistency is less immediately critical than *generalization*)...

Consistency

We say that an algorithm is consistent if

$$\forall \epsilon > 0 \quad \lim_{n \to \infty} \mathbb{P}\{I[f_n] - I[f_*] \ge \epsilon\} = 0$$

Universal Consistency

We say that an algorithm is universally consistent if for all probability p,

$$\forall \epsilon > 0 \quad \lim_{n \to \infty} \mathbb{P}\{I[f_n] - I[f_*] \ge \epsilon\} = 0$$

イロト イポト イヨト イヨト

The above requirements are asymptotic.

Error Rates

A more practical question is, how fast does the error decay? This can be expressed as

$$\mathbb{P}\{I[f_n] - I[f_*]\} \le \epsilon(n, \delta)\} \ge 1 - \delta.$$

Sample Complexity

Or equivalently, 'how many point do we need to achieve an error ϵ with a prescribed probability δ ?' This can expressed as

$$\mathbb{P}\{I[f_n] - I[f_*] \le \epsilon\} \ge 1 - \delta,$$

for $n = n(\epsilon, \delta)$.

ヘロト ヘ戸ト ヘヨト ヘヨト

The above requirements are asymptotic.

Error Rates

A more practical question is, how fast does the error decay? This can be expressed as

$\mathbb{P}\{I[f_n] - I[f_*]\} \le \epsilon(n, \delta)\} \ge 1 - \delta.$

Sample Complexity

Or equivalently, 'how many point do we need to achieve an error ϵ with a prescribed probability δ ?' This can expressed as

$$\mathbb{P}\{I[f_n] - I[f_*] \le \epsilon\} \ge 1 - \delta,$$

for $n = n(\epsilon, \delta)$.

ヘロト 人間 ト 人 ヨ ト 人 ヨ ト

The above requirements are asymptotic.

Error Rates

A more practical question is, how fast does the error decay? This can be expressed as

$$\mathbb{P}\{I[f_n] - I[f_*]\} \le \epsilon(n, \delta)\} \ge 1 - \delta.$$

Sample Complexity

Or equivalently, 'how many point do we need to achieve an error ϵ with a prescribed probability δ ?' This can expressed as

$$\mathbb{P}\{I[f_n] - I[f_*] \le \epsilon\} \ge 1 - \delta,$$

for $n = n(\epsilon, \delta)$.

・ 同 ト ・ ヨ ト ・ ヨ ト

The above requirements are asymptotic.

Error Rates

A more practical question is, how fast does the error decay? This can be expressed as

$$\mathbb{P}\{I[f_n] - I[f_*]\} \le \epsilon(n, \delta)\} \ge 1 - \delta.$$

Sample Complexity

Or equivalently, 'how many point do we need to achieve an error ϵ with a prescribed probability δ ?'

This can expressed as

$\mathbb{P}\{I[f_n] - I[f_*] \le \epsilon\} \ge 1 - \delta,$

for $n = n(\epsilon, \delta)$.

・ロ・ ・ 四・ ・ ヨ・ ・ ヨ・

э

The above requirements are asymptotic.

Error Rates

A more practical question is, how fast does the error decay? This can be expressed as

$$\mathbb{P}\{I[f_n] - I[f_*]\} \le \epsilon(n, \delta)\} \ge 1 - \delta.$$

Sample Complexity

Or equivalently, 'how many point do we need to achieve an error ϵ with a prescribed probability δ ?' This can expressed as

$$\mathbb{P}\{I[f_n] - I[f_*] \le \epsilon\} \ge 1 - \delta,$$

for $n = n(\epsilon, \delta)$.

ヘロト ヘワト ヘビト ヘビト
Empirical risk and Generalization

How do we design learning algorithms that work? One of the most natural ideas is ERM...

Empirical Risk

The empirical risk is a natural proxy (how good?) for the expected risk

$$I_n[f] = \frac{1}{n} \sum_{i=1}^n V(f(x_i), y_i).$$

Generalization Error

How good a proxy is captured by the generalization error,

$$\mathbb{P}\{|I[f_n] - I_n[f_n]| \le \epsilon\} \ge 1 - \delta,$$

for $n = n(\epsilon, \delta)$.

Empirical risk and Generalization

How do we design learning algorithms that work? One of the most natural ideas is ERM...

Empirical Risk

The empirical risk is a natural proxy (how good?) for the expected risk

$$I_n[f] = \frac{1}{n} \sum_{i=1}^n V(f(x_i), y_i).$$

Generalization Error

How good a proxy is captured by the generalization error,

$$\mathbb{P}\{|I[f_n] - I_n[f_n]| \le \epsilon\} \ge 1 - \delta,$$

for $n = n(\epsilon, \delta)$.

Empirical risk and Generalization

How do we design learning algorithms that work? One of the most natural ideas is ERM...

Empirical Risk

The empirical risk is a natural proxy (how good?) for the expected risk

$$I_n[f] = \frac{1}{n} \sum_{i=1}^n V(f(x_i), y_i).$$

Generalization Error

How good a proxy is captured by the generalization error,

$$\mathbb{P}\{|I[f_n] - I_n[f_n]| \le \epsilon\} \ge 1 - \delta,$$

for $n = n(\epsilon, \delta)$.

Some (Theoretical and Practical) Questions

• How do we go from here to an actual class of algorithms?

- Is minimizing the empirical error error on the data a good idea?
- Under which conditions is the empirical error a good proxy for the expected error?

ヘロト ヘアト ヘビト ヘビト

Some (Theoretical and Practical) Questions

- How do we go from here to an actual class of algorithms?
- Is minimizing the empirical error error on the data a good idea?
- Under which conditions is the empirical error a good proxy for the expected error?

ヘロト 人間 ト ヘヨト ヘヨト

Some (Theoretical and Practical) Questions

- How do we go from here to an actual class of algorithms?
- Is minimizing the empirical error error on the data a good idea?
- Under which conditions is the empirical error a good proxy for the expected error?

・ 同 ト ・ ヨ ト ・ ヨ ト …

- Part I: Basic Concepts and Notation
- Part II: Foundational Results
- Part III: Algorithms

・ 回 ト ・ ヨ ト ・ ヨ ト

ъ

Since classical statistics worries so much about consistency let us start here even if I do not think it is a practically important concept. Can we learn consistently any problem? Or equivalently do universally consistent algorithms exist? YES! Neareast neighbors, Histogram rules, SVM with (so called) universal kernels...

No Free Lunch Theorem

Given a number of points (and a confidence), can we always achieve a prescribed error? NO!

The last statement can be interpreted as follows: inference from finite samples can effectively performed if and only if the problem satisfies some a priori condition.

Since classical statistics worries so much about consistency let us start here even if I do not think it is a practically important concept. Can we learn consistently any problem? Or equivalently do universally consistent algorithms exist? YES! Neareast neighbors, Histogram rules, SVM with (so called) universal kernels...

No Free Lunch Theorem

Given a number of points (and a confidence), can we always achieve a prescribed error? NO!

The last statement can be interpreted as follows: inference from finite samples can effectively performed if and only if the problem satisfies some a priori condition.

Since classical statistics worries so much about consistency let us start here even if I do not think it is a practically important concept. Can we learn consistently any problem? Or equivalently do universally consistent algorithms exist? YES! Neareast neighbors, Histogram rules, SVM with (so called) universal kernels...

No Free Lunch Theorem

Given a number of points (and a confidence), can we always achieve a prescribed error?

NO!

The last statement can be interpreted as follows: inference from finite samples can effectively performed if and only if the problem satisfies some a priori condition.

Since classical statistics worries so much about consistency let us start here even if I do not think it is a practically important concept. Can we learn consistently any problem? Or equivalently do universally consistent algorithms exist? YES! Neareast neighbors, Histogram rules, SVM with (so called) universal kernels...

No Free Lunch Theorem

Given a number of points (and a confidence), can we always achieve a prescribed error? NO!

The last statement can be interpreted as follows: inference from finite samples can effectively performed if and only if the problem satisfies some a priori condition.

In many learning algorithms (not all!) we need to choose a suitable space of hypotheses \mathcal{H} .

The **hypothesis space** \mathcal{H} is the space of functions that we allow our algorithm to "look at". For many algorithms (such as optimization algorithms) it is the space the algorithm is allowed to search. As we will see in future classes, it is often important to choose the hypothesis space as a function of the amount of data *n* available.

・ 同 ト ・ ヨ ト ・ ヨ ト

In many learning algorithms (not all!) we need to choose a suitable space of hypotheses \mathcal{H} .

The **hypothesis space** \mathcal{H} is the space of functions that we allow our algorithm to "look at". For many algorithms (such as optimization algorithms) it is the space the algorithm is allowed to search. As we will see in future classes, it is often important to choose the hypothesis space as a function of the amount of data *n* available.

Examples: linear functions, polynomial, RBFs, Sobolev Spaces...

_earning algorithm

A learning algorithm A is then a map from the data space to \mathcal{H} ,

 $A(S_n) = f_n \in \mathcal{H}.$

Tomaso Poggio The Learning Problem and Regularization

・ロト ・ 理 ト ・ ヨ ト ・

3

Examples: linear functions, polynomial, RBFs, Sobolev Spaces...

Learning algorithm

A learning algorithm A is then a map from the data space to \mathcal{H} ,

$$A(S_n) = f_n \in \mathcal{H}.$$

くロト (過) (目) (日)

æ

ERM

A prototype algorithm in statistical learning theory is Empirical Risk Minimization:

 $\min_{f\in\mathcal{H}}I_n[f].$

How do we choose \mathcal{H} ? How do we design A?

くロト (過) (目) (日)

ъ

Given a function f, a loss function V, and a probability distribution μ over Z, the **expected or true error** of f is:

$$I[f] = \mathbb{E}_{z}V[f,z] = \int_{Z}V(f,z)d\mu(z)$$

which is the **expected loss** on a new example drawn at random from μ .

We would like to make I[f] small, but in general we do not know μ .

Given a function f, a loss function V, and a training set S consisting of n data points, the **empirical error** of f is:

$$I_{\mathcal{S}}[f] = \frac{1}{n} \sum V(f, z_i)$$

・ 同 ト ・ ヨ ト ・ ヨ ト ・

Reminder: Generalization

A natural requirement for f_S is distribution independent **generalization**

 $\lim_{n\to\infty} |I_{\mathcal{S}}[f_{\mathcal{S}}] - I[f_{\mathcal{S}}]| = 0 \text{ in probability}$

This is equivalent to saying that for each *n* there exists a ε_n and a $\delta(\varepsilon)$ such that

$$\mathbb{P}\left\{\left|I_{S_n}[f_{S_n}] - I[f_{S_n}]\right| \ge \varepsilon_n\right\} \le \delta(\varepsilon_n),\tag{1}$$

with ε_n and δ going to zero for $n \to \infty$.

In other words, the training error for the solution must converge to the expected error and thus be a "proxy" for it. Otherwise the solution would not be "predictive".

A desirable additional requirement is **consistency**

$$\varepsilon > 0 \lim_{n \to \infty} \mathbb{P}\left\{ I[f_{\mathcal{S}}] - \inf_{f \in \mathcal{H}} I[f] \ge \varepsilon \right\} = 0.$$

In addition to the key property of generalization, a "good" learning algorithm should also be *stable*: f_S should depend continuously on the training set *S*. In particular, changing one of the training points should affect less and less the solution as *n* goes to infinity. Stability is a good requirement for the learning problem and, in fact, for any mathematical problem. We open here a small parenthesis on stability and well-posedness.

・ 同 ト ・ ヨ ト ・ ヨ ト

A problem is **well-posed** if its solution:

- exists
- is unique
- depends continuously on the data (e.g. it is stable)

A problem is **ill-posed** if it is not well-posed. In the context of this class, well-posedness is mainly used to mean *stability* of the solution.

・ 回 ト ・ ヨ ト ・ ヨ ト

More on well-posed and ill-posed problems

Hadamard introduced the definition of ill-posedness. Ill-posed problems are typically inverse problems.

As an example, assume g is a function in Y and u is a function in X, with Y and X Hilbert spaces. Then given the linear, continuous operator L, consider the equation

$$g = Lu$$
.

The direct problem is is to compute g given u; the inverse problem is to compute u given the data g. In the learning case L is somewhat similar to a "sampling" operation and the inverse problem becomes the problem of finding a function that takes the values

$$f(x_i) = y_i, i = 1, ..., n$$

The inverse problem of finding *u* is well-posed when

- the solution exists,
- is unique and
- is stable, that is depends continuously on the initial data g

Given a training set *S* and a function space \mathcal{H} , empirical risk minimization as we have seen is the class of algorithms that look at *S* and select f_S as

$$f_{\mathcal{S}} = \arg\min_{f\in\mathcal{H}} I_{\mathcal{S}}[f].$$

For example linear regression is ERM when $V(z) = (f(x) - y)^2$ and *H* is space of linear functions f = ax.

・ 同 ト ・ ヨ ト ・ ヨ ト ・

For ERM to represent a "good" class of learning algorithms, the solution should

- generalize
- exist, be unique and especially be stable (well-posedness), according to some definition of stability.

(同) くほり くほう

ERM and generalization: given a certain number of samples...



Tomaso Poggio The Learning Problem and Regularization

...suppose this is the "true" solution...



프 > 프

... but suppose ERM gives this solution.



Tomaso Poggio The Learning Problem and Regularization

프 🕨 🗆 프

Under which conditions the ERM solution converges with increasing number of examples to the true solution? In other words...what are the conditions for generalization of ERM?



ERM and stability: given 10 samples...



Tomaso Poggio The Learning Problem and Regularization

æ

э

...we can find the smoothest interpolating polynomial (which degree?).



Tomaso Poggio The Learning Problem and Regularization

But if we perturb the points slightly...



Tomaso Poggio The Learning Problem and Regularization

æ

э

...the solution changes a lot!



Tomaso Poggio The Learning Problem and Regularization

▶ < Ξ >

æ

If we restrict ourselves to degree two polynomials...



Tomaso Poggio The Learning Problem and Regularization

ъ

э

...the solution varies only a small amount under a small perturbation.



Tomaso Poggio The Learning Problem and Regularization

Since Tikhonov, it is well-known that a generally ill-posed problem such as ERM, can be guaranteed to be well-posed and therefore *stable* by an appropriate choice of \mathcal{H} . For example, compactness of \mathcal{H} guarantees stability. It seems intriguing that Vapnik's (see also Cucker and Smale) *classical conditions for consistency of ERM* – thus quite a different property – consist of appropriately restricting \mathcal{H} . It seems that the same restrictions that make the approximation of the data stable, may provide solutions that generalize...

・ 回 ト ・ ヨ ト ・ ヨ ト

We would like to have a hypothesis space that yields generalization. Loosely speaking this would be a *H* for which the solution of ERM, say f_S is such that $|I_S[f_S] - I[f_S]|$ converges to zero in probability for *n* increasing. Note that the above requirement is NOT the law of large numbers; the requirement for a fixed *f* that $|I_S[f] - I[f]|$ converges to zero in probability for *n* increasing IS the law of large numbers.

ヘロト ヘ戸ト ヘヨト ヘヨト

ERM: conditions for well-posedness (stability) and predictivity (generalization) in the case of regression and classification

- The theorem (Vapnik et al.) says that a proper choice of the hypothesis space H ensures generalization of ERM (and consistency since for ERM generalization is necessary and sufficient for consistency and viceversa). Other results characterize uGC classes in terms of measures of complexity or capacity of H (such as VC dimension).
- A separate theorem (Niyogi, Mukherjee, Rifkin, Poggio) says that stability (defined in a specific way) of (supervised) ERM is sufficient and necessary for generalization of ERM. Thus with the appropriate definition of stability, stability and generalization are equivalent for ERM; stability and H uGC are also equivalent.

Thus the two desirable conditions for a supervised learning algorithm – generalization and stability – are equivalent (and they correspond to the same constraints on \mathcal{H}).

イロト イポト イヨト イヨト
Key Theorem(s) Illustrated



Tomaso Poggio The Learning Problem and Regularization

ヘロト 人間 とくほとくほとう



・ロト・「日下・「日下・「日下」 シック

L

The "equivalence" between generalization and stability gives us a an approach to predictive algorithms. It is enough to remember that regularization is the classical way to restore well posedness. Thus regularization becomes a way to ensure generalization. Regularization in general means retricting H, as we have in fact done for ERM. There are two standard approaches in the field of ill-posed problems that ensure for ERM *well-posedness* (and *generalization*) by constraining the hypothesis space \mathcal{H} . The direct way – minimize the empirical error subject to f in a ball in an appropriate \mathcal{H} – is called *lvanov* regularization. The indirect way is Tikhonov regularization (which is not strictly ERM).

ヘロト ヘ戸ト ヘヨト ヘヨト

Ivanov and Tikhonov Regularization

ERM finds the function in (\mathcal{H}) which minimizes

$$\frac{1}{n}\sum_{i=1}^{n}V(f(x_i),y_i)$$

which in general – for arbitrary hypothesis space H – is *ill-posed*.

Ivanov regularizes by finding the function that minimizes

$$\frac{1}{n}\sum_{i=1}^{n}V(f(x_i), y_i)$$

while satisfying $\mathcal{R}(f) \leq A$.

 Tikhonov regularization minimizes over the hypothesis space H, for a fixed positive parameter γ, the regularized functional

$$\frac{1}{n}\sum_{i=1}^{n}V(f(x_i), y_i) + \gamma \mathcal{R}(f).$$
(2)

イロン 不得 とくほ とくほう 一日

 $\mathcal{R}(f)$ is the regulirizer, a penalization on f. In this course we will mainly discuss the case $\mathcal{R}(f) = \|f\|_{K}^{2}$ where $\|f\|_{K}^{2}$ is the norm in the Reproducing Kernel Hilbert Space (RKHS) \mathcal{H} , defined by the kernel K.

As we will see in future classes

- Tikhonov regularization ensures well-posedness eg existence, uniqueness and especially *stability* (in a very strong form) of the solution
- Tikhonov regularization ensures generalization
- Tikhonov regularization is closely related to but different from – Ivanov regularization, eg ERM on a hypothesis space H which is a ball in a RKHS.

伺き くほき くほう

Intelligent behavior (at least learning) consists of optimizing under constraints. Constraints are key for solving computational problems; constraints are key for prediction. Constraints may correspond to rather general symmetry properties of the problem (eg time invariance, space invariance, invariance to physical units (pai theorem), universality of numbers and metrics implying normalization, etc.)

- Key questions at the core of learning theory:
 - generalization and predictivity not explanation
 - probabilities are unknown, only data are given
 - which constraints are needed to ensure generalization (therefore which hypotheses spaces)?
 - regularization techniques result usually in *computationally "nice"* and *well-posed* optimization problems

ヘロト 人間 ト くほ ト くほ トー

Unlike statistical learning theory the Bayesian approach does not emphasize

- the issue of generalization (following the tradition in statistics of explanatory statistics);
- that probabilities are not known and that only data are known: assuming a specific distribution is a very strong – unconstrained by any Bayesian theory – seat-of-the-pants guess;
- the question of which priors are needed to ensure generalization;
- that the resulting optimization problems are often *computationally intractable* and possibly ill-posed optimization problems (for instance not unique).

・ 同 ト ・ ヨ ト ・ ヨ ト …

- Part I: Basic Concepts and Notation
- Part II: Foundational Results
- Part III: Algorithms

INSTEAD....

・ 同 ト ・ ヨ ト ・ ヨ ト

3

In addition to the hypothesis space \mathcal{H} , the space we allow our algorithms to search, we define...

The **target space** \mathcal{T} is a space of functions, chosen a priori in any given problem, that is assumed to contain the "true" function f_0 that minimizes the risk. Often, \mathcal{T} is chosen to be all functions in L_2 , or all differentiable functions. Notice that the "true" function if it exists is defined by $\mu(z)$, which contains all the relevant information.

・ 同 ト ・ ヨ ト ・ ヨ ト

Let $f_{\mathcal{H}}$ be the function in \mathcal{H} with the smallest true risk. We have defined the **generalization error** to be $I_S[f_S] - I[f_S]$. We define the **sample error** to be $I[f_S] - I[f_{\mathcal{H}}]$, the difference in true risk between the best function in \mathcal{H} and the function in \mathcal{H} we actually find. This is what we pay because our finite sample does not give us enough information to choose to the "best" function in \mathcal{H} . We'd like this to be small. *Consistency* – defined earlier – is equivalent to the sample error going to zero for $n \to \infty$.

A main goal in classical learning theory (Vapnik, Smale, ...) is "bounding" the generalization error. Another goal – for learning theory and statistics – is bounding the sample error, that is determining conditions under which we can state that $I[f_S] - I[f_H]$ will be small (with high probability).

As a simple rule, we expect that if \mathcal{H} is "well-behaved", then, as *n* gets large the sample error will become small.

ヘロア ヘビア ヘビア・

Let f_0 be the function in \mathcal{T} with the smallest true risk. We define the **approximation error** to be $I[f_{\mathcal{H}}] - I[f_0]$, the difference in true risk between the best function in \mathcal{H} and the best function in \mathcal{T} . This is what we pay when \mathcal{H} is smaller than \mathcal{T} . We'd like this error to be small too. In much of the following we can assume that $I[f_0] = 0$.

We will focus less on the approximation error in 9.520, but we will explore it.

As a simple rule, we expect that as \mathcal{H} grows bigger, the approximation error gets smaller. If $\mathcal{T} \subseteq \mathcal{H}$ – which is a situation called *the realizable setting* –the approximation error is zero.

ヘロア ヘビア ヘビア・

We define the **error** to be $I[f_S] - I[f_0]$, the difference in true risk between the function we actually find and the best function in \mathcal{T} . We'd really like this to be small. As we mentioned, often we can assume that the **error** is simply $I[f_S]$.

The error is the sum of the sample error and the approximation error:

$$I[f_{S}] - I[f_{0}] = (I[f_{S}] - I[f_{H}]) + (I[f_{H}] - I[f_{0}])$$

If we can make both the approximation and the sample error small, the error will be small. There is a tradeoff between the approximation error and the sample error...

・ 回 ト ・ ヨ ト ・ ヨ ト

It should already be intuitively clear that making \mathcal{H} big makes the approximation error small. This implies that we can (help) make the error small by making \mathcal{H} big.

On the other hand, we will show that making \mathcal{H} small will make the sample error small. In particular for ERM, if \mathcal{H} is a uGC class, the generalization error and the sample error will go to zero as $n \to \infty$, but how quickly depends directly on the "size" of \mathcal{H} . This implies that we want to keep \mathcal{H} as small as possible. (Furthermore, \mathcal{T} itself may or may not be a uGC class.) Ideally, we would like to find the optimal tradeoff between these conflicting requirements.

ヘロト ヘ戸ト ヘヨト ヘヨト

Generalization error is $I_S[f_S] - I[f_S]$. Sample error is $I[f_S] - I[f_H]$ Approximation error is $I[f_H] - I[f_0]$ Error is $I[f_S] - I[f_0] = (I[f_S] - I[f_H]) + (I[f_H] - I[f_0])$

伺 とくきとくきと

- Part I: Basic Concepts and Notation
- Part II: Foundational Results
- Part III: Algorithms

・ 同 ト ・ ヨ ト ・ ヨ ト

ъ

We are going to look at hypotheses spaces which are reproducing kernel Hilbert spaces.

- RKHS are **Hilbert spaces** of **point-wise defined** functions.
- They can be defined via a **reproducing kernel**, which is a symmetric positive definite function.

$$\sum_{i,j=1}^n c_i c_j K(t_i, t_j) \ge 0$$

for any $n \in \mathbb{N}$ and choice of $t_1, ..., t_n \in X$ and $c_1, ..., c_n \in \mathbb{R}$. • functions in the space are (the completion of) linear combinations

$$f(x) = \sum_{i=1}^{p} K(x, x_i) c_i.$$

• the norm in the space is a natural measure of complexity

< 回 > < 回 > < 回 > … 回

We are going to look at hypotheses spaces which are reproducing kernel Hilbert spaces.

- RKHS are **Hilbert spaces** of **point-wise defined** functions.
- They can be defined via a **reproducing kernel**, which is a symmetric positive definite function.

$$\sum_{i,j=1}^n c_i c_j K(t_i,t_j) \ge 0$$

for any $n \in \mathbb{N}$ and choice of $t_1, ..., t_n \in X$ and $c_1, ..., c_n \in \mathbb{R}$.

 functions in the space are (the completion of) linear combinations

$$f(x) = \sum_{i=1}^{p} K(x, x_i) c_i.$$

• the norm in the space is a natural measure of complexity

< 回 > < 回 > < 回 > … 回

We are going to look at hypotheses spaces which are reproducing kernel Hilbert spaces.

- RKHS are **Hilbert spaces** of **point-wise defined** functions.
- They can be defined via a **reproducing kernel**, which is a symmetric positive definite function.

$$\sum_{i,j=1}^n c_i c_j K(t_i,t_j) \ge 0$$

for any $n \in \mathbb{N}$ and choice of $t_1, ..., t_n \in X$ and $c_1, ..., c_n \in \mathbb{R}$.

functions in the space are (the completion of) linear combinations

$$f(x) = \sum_{i=1}^{p} K(x, x_i) c_i.$$

• the norm in the space is a natural measure of complexity

We are going to look at hypotheses spaces which are reproducing kernel Hilbert spaces.

- RKHS are **Hilbert spaces** of **point-wise defined** functions.
- They can be defined via a **reproducing kernel**, which is a symmetric positive definite function.

$$\sum_{i,j=1}^n c_i c_j K(t_i,t_j) \ge 0$$

for any $n \in \mathbb{N}$ and choice of $t_1, ..., t_n \in X$ and $c_1, ..., c_n \in \mathbb{R}$.

functions in the space are (the completion of) linear combinations

$$f(x) = \sum_{i=1}^{p} K(x, x_i) c_i.$$

the norm in the space is a natural measure of complexity

Very common examples of symmetric pd kernels are • Linear kernel

$$K(x,x')=x\cdot x'$$

Gaussian kernel

$$\mathcal{K}(\mathbf{x},\mathbf{x}') = \mathrm{e}^{-rac{\|\mathbf{x}-\mathbf{x}'\|^2}{\sigma^2}}, \qquad \sigma > \mathbf{0}$$

Polynomial kernel

$$K(x,x') = (x \cdot x' + 1)^d, \qquad d \in \mathbb{N}$$

For specific applications, designing an effective kernel is a challenging problem.

< 回 > < 回 > < 回 > .

Often times kernels, are defined through a dictionary of features

$$\mathcal{D} = \{\phi_j, i = 1, \dots, p \mid \phi_j : X \to \mathbb{R}, \forall j\}$$

setting

$$\mathcal{K}(\mathbf{x},\mathbf{x}')=\sum_{i=1}^p\phi_j(\mathbf{x})\phi_j(\mathbf{x}').$$

・ 同 ト ・ ヨ ト ・ ヨ ト …

3

We can regularize by explicitly restricting the hypotheses space \mathcal{H} — for example to a ball of radius R.

Ivanov regularization

$$\min_{f\in\mathcal{H}}\frac{1}{n}\sum_{i=1}^{n}V(f(x_i),y_i)$$

subject to

$$\|f\|_{\mathcal{H}}^2 \leq \mathbf{R}.$$

The above algorithm corresponds to a constrained optimization problem.

・ロト ・聞 ト ・ ヨト ・ ヨトー

Regularization can also be done implicitly via penalization

Tikhonov regularizarion

$$\arg\min_{f\in\mathcal{H}}\frac{1}{n}\sum_{i=1}^{n}V(f(x_i),y_i)+\frac{\lambda}{\|f\|_{\mathcal{H}}^2}.$$

 $\boldsymbol{\lambda}$ is the regularization parameter trading-off between the two terms.

The above algorithm can be seen as the Lagrangian formulation of a constrained optimization problem.

・ 同 ト ・ ヨ ト ・ ヨ ト

An important result

The minimizer over the RKHS \mathcal{H} , f_S , of the regularized empirical functional

 $I_{\mathcal{S}}[f] + \lambda \|f\|_{\mathcal{H}}^2,$

can be represented by the expression

$$f_n(x) = \sum_{i=1}^n c_i K(x_i, x),$$

for some $(c_1, \ldots, c_n) \in \mathbb{R}$.

Hence, minimizing over the (possibly infinite dimensional) Hilbert space, *boils down to minimizing over* \mathbb{R}^n .

ヘロト ヘ戸ト ヘヨト ヘヨト

The way the coefficients $\mathbf{c} = (c_1, \dots, c_n)$ are computed depend on the loss function choice.

• RLS: Let Let $\mathbf{y} = (y_1, \dots, y_n)$ and $\mathbf{K}_{i,j} = K(x_i, x_j)$ then $\mathbf{c} = (\mathbf{K} + \lambda n l)^{-1} \mathbf{y}$.

• SVM: Let $\alpha_i = y_i c_i$ and $\mathbf{Q}_{i,j} = y_i K(x_i, x_j) y_j$

$$\max_{\alpha \in \mathbb{R}^{n}} \sum_{i=1}^{n} \alpha_{i} - \frac{1}{2} \alpha^{T} \mathbf{Q} \alpha$$

subject to :
$$\sum_{i=1}^{n} y_{i} \alpha_{i} = 0$$
$$0 \le \alpha_{i} \le C \qquad i = 1, \dots, n$$

◆□▶ ◆□▶ ★ □▶ ★ □▶ → □ → の Q ()

The way the coefficients $\mathbf{c} = (c_1, \dots, c_n)$ are computed depend on the loss function choice.

- RLS: Let Let $\mathbf{y} = (y_1, \dots, y_n)$ and $\mathbf{K}_{i,j} = \mathcal{K}(x_i, x_j)$ then $\mathbf{c} = (\mathbf{K} + \lambda n l)^{-1} \mathbf{y}$.
- SVM: Let $\alpha_i = y_i c_i$ and $\mathbf{Q}_{i,j} = y_i K(x_i, x_j) y_j$

$$\max_{\alpha \in \mathbb{R}^{n}} \sum_{i=1}^{n} \alpha_{i} - \frac{1}{2} \alpha^{T} \mathbf{Q} \alpha$$

subject to :
$$\sum_{i=1}^{n} y_{i} \alpha_{i} = 0$$
$$0 \le \alpha_{i} \le C \qquad i = 1, \dots, n$$

◆□▶ ◆□▶ ★ □▶ ★ □▶ → □ → の Q ()

Empirical risk minimization is ML.

$$p(\mathbf{Y}|f, \mathbf{X}) \propto e^{-\frac{1}{2}\sum_{i=1}^{N}(y_i - f(x_i))^2}$$

Linear RLS is MAP.

$$p(\mathbf{Y}, f | \mathbf{X}) \propto e^{-\frac{1}{2} \sum_{i=1}^{N} (y_i - \langle x_i, \theta \rangle)^2} \cdot e^{-\frac{\lambda}{2} \theta^T \theta}$$

Kernel RLS is also MAP.

$$p(\mathbf{Y}, f | \mathbf{X}) \propto e^{-\frac{1}{2} \sum_{i=1}^{N} (y_i - f(x_i)^2} \cdot e^{-\frac{\lambda}{2} \|f\|_{\mathcal{H}}^2}$$

More generally we can consider:

 $I_n(f) + \lambda R(f)$

where, R(f) is a regularizing functional.

- Sparsity based methods
- Manifold learning
- Multiclass
- ...

・ 同 ト ・ ヨ ト ・ ヨ ト …

æ

- statistical learning as a foundational framework to predict from data
- a proxy for predictivity is the empirical error iff generalization holds for the class of algorithms
- stability and generalization are equivalent
- regularization as a fundamental tool in learning algorithm to ensure stability and generalization

・ 同 ト ・ ヨ ト ・ ヨ ト ・

Generalization error is $I_S[f_S] - I[f_S]$. Sample error is $I[f_S] - I[f_H]$ Approximation error is $I[f_H] - I[f_0]$ Error is $I[f_S] - I[f_0] = (I[f_S] - I[f_H]) + (I[f_H] - I[f_0])$

伺 とくきとくきと

Final (optional) Remarks

Tomaso Poggio The Learning Problem and Regularization

ヘロン 人間 とくほ とくほ とう

æ

Intelligent behavior (at least learning) consists of optimizing under constraints. Constraints are key for solving computational problems; constraints are key for prediction. Constraints may correspond to rather general symmetry properties of the problem (eg time invariance, space invariance, invariance to physical units (π theorem), universality of numbers and metrics implying normalization, etc.)

・ 同 ト ・ ヨ ト ・ ヨ ト …

ERM: conditions for well-posedness (stability) and predictivity (generalization) in the case of regression and classification

Theorem [Vapnik and Červonenkis (71), Alon et al (97), Dudley, Giné, and Zinn (91)]

if

A (necessary) and sufficient condition for generalization (and consistency) of ERM is that \mathcal{H} is uGC. **Definition** \mathcal{H} is a (weak) uniform Glivenko-Cantelli (uGC) class

$$\forall \varepsilon > 0 \lim_{n \to \infty} \sup_{\mu} \mathbb{P}_{\mathcal{S}} \left\{ \sup_{f \in \mathcal{H}} |I[f] - I_{\mathcal{S}}[f]| > \varepsilon \right\} = 0.$$

ヘロト 人間 ト くほ ト くほ トー

Uniform Glivenko-Cantelli Classes

We say that \mathcal{H} is a uniform Glivenko-Cantelli (uGC) class, if for all p,

$$\forall \epsilon > 0 \lim_{n \to \infty} \mathbb{P} \left\{ \sup_{f \in \mathcal{H}} |I[f] - I_n[f]| > \epsilon \right\} = 0.$$

A necessary and sufficient condition for consistency of ERM is that \mathcal{H} is uGC. See: [Vapnik and Červonenkis (71), Alon et al (97), Dudley, Giné, and Zinn (91)].

In turns the UGC property is equivalent to requiring \mathcal{H} to have finite capacity: V_{γ} dimension in general and VC dimension in classification.

イロト イポト イヨト イヨト

Uniform Glivenko-Cantelli Classes

We say that \mathcal{H} is a uniform Glivenko-Cantelli (uGC) class, if for all p,

$$\forall \epsilon > 0 \lim_{n \to \infty} \mathbb{P} \left\{ \sup_{f \in \mathcal{H}} |I[f] - I_n[f]| > \epsilon \right\} = 0.$$

A necessary and sufficient condition for consistency of ERM is that \mathcal{H} is uGC. See: [Vapnik and Červonenkis (71), Alon et al (97), Dudley, Giné, and Zinn (91)].

In turns the UGC property is equivalent to requiring \mathcal{H} to have finite capacity: V_{γ} dimension in general and VC dimension in classification.

notation: *S* training set, $S^{i,z}$ training set obtained replacing the *i*-th example in *S* with a new point z = (x, y).

Definition

We say that an algorithm \mathcal{A} has **uniform stability** β (is β -stable) if

$$\forall (S,z) \in \mathcal{Z}^{n+1}, \ \forall i, \ \sup_{z' \in Z} |V(f_S,z') - V(f_{S^{i,z}},z')| \leq \beta$$

ヘロト ヘアト ヘビト ヘビト

æ
$$z = (x, y)$$

 $S = z_1, ..., z_n$
 $S^i = z_1, ..., z_{i-1}, z_{i+1}, ..., z_n$

CV Stability

A learning algorithm *A* is CV_{loo} stable if for each *n* there exists a $\beta_{CV}^{(n)}$ and a $\delta_{CV}^{(n)}$ such that for all *p*

$$\mathbb{P}\left\{|m{V}(f_{\mathcal{S}^{i}}, z_{i}) - m{V}(f_{\mathcal{S}}, z_{i})| \leq eta_{CV}^{(n)}
ight\} \geq 1 - \delta_{CV}^{(n)},$$

with $\beta_{CV}^{(n)}$ and $\delta_{CV}^{(n)}$ going to zero for $n \to \infty$.

・ 同 ト ・ ヨ ト ・ ヨ ト ・

In the above reasoning the kernel and the hypotheses space define a representation/parameterization of the problem and hence play a special role.

Where do they come from?

- There are a few off the shelf choices (Gaussian, polynomial etc.)
- Often they are the product of problem specific engineering.

Are there principles— applicable in a wide range of situations to design effective data representation?

・ 同 ト ・ 三 ト ・

In the above reasoning the kernel and the hypotheses space define a representation/parameterization of the problem and hence play a special role.

Where do they come from?

- There are a few off the shelf choices (Gaussian, polynomial etc.)
- Often they are the product of problem specific engineering.

Are there principles— applicable in a wide range of situations to design effective data representation?

・ 同 ト ・ ヨ ト ・ ヨ

In the above reasoning the kernel and the hypotheses space define a representation/parameterization of the problem and hence play a special role.

Where do they come from?

- There are a few off the shelf choices (Gaussian, polynomial etc.)
- Often they are the product of problem specific engineering.

Are there principles— applicable in a wide range of situations to design effective data representation?

In the above reasoning the kernel and the hypotheses space define a representation/parameterization of the problem and hence play a special role.

Where do they come from?

- There are a few off the shelf choices (Gaussian, polynomial etc.)
- Often they are the product of problem specific engineering.

Are there principles– applicable in a wide range of situations– to design effective data representation?