# Administrative Information

HW #1 is due tonight, by midnight. Please submit via Gradescope. The course access code is on HW #1 and Slack

The project description is now up. Feel free to email me if you'd like to chat and/or run something by me

The first deadline is April 9th, a one page project proposal

# Stability, Regularization and Double Descent

Today we will study generalization thru the following framework:

"If the hypothesis I learn does not depend on any one data point too much, then it must generalize"

Later we will see how stability comes about through explicit / implicit regularization

As usual we will be interested in the risk:

$$R(w) = \underset{z=(x,y) \sim D}{\mathbb{E}} \left[ \ell(z, w) \right]$$

and the empirical risk:

$$R_s(w) = \frac{1}{N} \sum_i \ell(z_i, w)$$

and the generalization error:

$$\varepsilon_{gen}(w) \stackrel{\Delta}{=} R(w) - R_s(w)$$

Or alternatively $R(w) = R_S(w) + \mathcal{E}_{gen}(w)$

First, we will show that generalization error is an average-case notion of stability

Consider an algorithm $A$ that we think of as a function

$$A: (X \times Y)^N \to \Omega$$

set of allowable params.

Now, similar to what we did in the context of Rademacher complexity, can create <u>hybrid samples</u>

$$S = (z_1, \dots, z_N)$$

$$S' = (z_1', \dots, z_N')$$ aka "ghost samples"

Now suppose we swap in the $i^{th}$ sample from $S'$:

$$S^{(i)} = (z_1, z_2, \ldots z_{i-1}, z_i', z_{i+1}, \ldots, z_N)$$

This leads us to:

<u>def</u>: The <u>average stability</u> of $A$ is

$$\Delta(A) = \mathop{\mathbb{E}}_{S, S'} \left[ \frac{1}{N} \sum_{i=1}^{N} \left( \ell(A(S), z_i') - \ell(A(S^{(i)}), z_i') \right) \right]$$

i.e. it is how much our loss changes on avg., when we swap out a single sample.

This is just generalization error in disguise.

<u>Theorem 1</u>   $\mathop{\mathbb{E}}_{S} \left[ \varepsilon_{gen}(A(S)) \right] = \Delta(A)$

<u>Proof</u>: Using the definition of generalization error, we have:

$$\mathop{\mathbb{E}}_{S} \left[ \varepsilon_{gen}(A(S)) \right] = \mathop{\mathbb{E}}_{S} \left[ R(A(S)) - R_S(A(S)) \right]$$

Expanding each term:

$$\mathbb{E}_S\left[R_S(A(S))\right] = \mathbb{E}_S\left[\frac{1}{N}\sum_i \ell(A(S), z_i)\right] \quad (*)$$

$$\mathbb{E}_S\left[R(A(S))\right] = \mathbb{E}_{S,S'}\left[\frac{1}{N}\sum_i \ell(A(S), z_i')\right] \quad (1)$$

The key point is because $z_i$ and $z_i'$ are identically distributed, we have

$$\mathbb{E}_S\left[\ell(A(S), z_i)\right] = \mathbb{E}_{S,S'}\left[\ell(A(S^{(i)}), z_i')\right]$$

Thus rewriting $(*)$ we have

$$\mathbb{E}_S\left[R_S(A(S))\right] = \mathbb{E}_{S,S'}\left[\frac{1}{N}\sum_i \ell(A(S^{(i)}), z_i')\right]$$

$$(2)$$

And combining (1) and (2), we're done

∎

It is usually challenging to work with average stability directly

Can instead use a _stronger_ notion

def: The uniform stability of $A$ is

$$\Delta_{unif}(A) = \sup_{S, S'} \sup_{i \in [N]} \left| \ell(A(S), z_i') - \ell(A(S^{(i)}), z_i') \right|$$

This is a worst-case notion of stability

i.e. the risk never changes too much, regardless of the samples you were given

Since this is a strengthening, we have:

Corollary 1: $\mathbb{E}_S \left[ \varepsilon_{gen}(A(S)) \right] \leq \Delta_{unif}(A)$

You can actually get bounds that hold with high probability:

Theorem 2 [Feldman, Vondrak] Let $\gamma = \Delta_{unif}(A)$

Then

$$\mathbb{P}_S \left[ \varepsilon_{gen}(A(S)) \geq \gamma \log N \log \frac{N}{\delta} + \sqrt{\frac{\log^{1/\delta}}{N}} \right] \leq \delta$$

where $S$ is an iid sample of size $N$

Earlier bounds of Bousquet and Elisseeff were only meaningful when $\gamma$ was somewhat large

In contrast, the bound above is best possible even when $\gamma = \frac{1}{N}$, i.e.

<span style="color:red">generalization error</span> $\sim$ <span style="color:red">sampling error</span>

## Stability of Empirical Risk Minimization

An important algorithm to analyze is

$$A : S \longrightarrow \underset{\omega}{\text{argmin}} \; R_S(\omega)$$

i.e. empirical risk minimization

As we discussed, we can't always solve this e.g. it is often nonconvex!

Next, we'll see that in some cases where we can, it is automatically stable

Following Shalev-Shwartz, Shamir, Srebro and Sridharan, we'll show:

Theorem 3 [Shalev-Shwartz et al]: Suppose the loss $\ell$ is $\alpha$-strongly convex and L-Lipschitz. Then the ERM

$$\hat{w}_S \overset{\Delta}{=} \underset{w \in \Omega}{\arg\min} \; \frac{1}{N} \sum_i \ell(w, z_i)$$

is uniformly stable, with $\Delta_{unif}(ERM) \leq \frac{4L^2}{\alpha N}$

Recall:

$\underline{\alpha\text{-strongly convex:}}$

$$\ell(w') \geq \ell(w) + \nabla \ell(w)^T (w'-w) + \frac{\alpha}{2} \| w'-w \|_2^2$$

$\underline{\text{L-Lipschitz:}}$

$$|\ell(w') - \ell(w)| \leq L \| w'-w \|_2$$

First we'll need an elementary fact:

## Fact 1:

If $f(w)$ and $g(w)$ are both $\alpha$-strongly convex then any convex combination

$$h(w) \stackrel{\Delta}{=} \lambda f(w) + (1-\lambda) g(w)$$

is also $\alpha$-strongly convex

**Proof:** This follows from the definition because

$$h(w') = \lambda f(w') + (1-\lambda) g(w')$$

$$\geq \lambda \left( f(w) + \nabla f(w)^T (w'-w) + \frac{\alpha}{2} \|w'-w\|_2^2 \right)$$

$$+ (1-\lambda) \left( g(w) + \nabla f(w)^T (w'-w) + \frac{\alpha}{2} \|w'-w\|_2^2 \right)$$

Collect terms
$$\underset{=}{} h(w) + \nabla h(w)^T (w'-w) + \frac{\alpha}{2} \|w'-w\|_2^2$$

$\blacksquare$

## Proof of Theorem 3:

First consider

$$R_S(\hat{w}_{S^{(i)}}) - R_S(\hat{w}_S)$$

Fact 1 tells us $R_S$ is $\alpha$-strongly convex,

$$R_S(\hat{w}_{S^{(i)}}) - R_S(\hat{w}_S) \geq \nabla R_S(w)\Big|_{\hat{w}_S}^T \overcancel{(\hat{w}_{S^{(i)}} - \hat{w}_S)}^{0}$$
$$+ \frac{\alpha}{2} \|\hat{w}_{S^{(i)}} - \hat{w}_S\|_2^2$$

From the definition of the ERM, we have:

$$\nabla R_S(w)\Big|_{\hat{w}_S} = 0$$

Thus we get

$$R_S(\hat{w}_{S^{(i)}}) - R(\hat{w}_S) \geq \frac{\alpha}{2} \|\hat{w}_{S^{(i)}} - \hat{w}_S\|_2^2 \quad (\text{*})$$

Next we'll upper bound the difference in risk. Let's pull out the $i^{th}$ term

$$R_S(\hat{w}_{S^{(i)}}) - R_S(\hat{w}_S) =$$

$$\frac{1}{N}\left(\ell(\hat{w}_{S^{(i)}}, z_i) - \ell(\hat{w}_S, z_i)\right) + \frac{1}{N}\sum_{j \neq i} \ell(\hat{w}_{S^{(i)}}, z_j)$$
$$- \ell(\hat{w}_S, z_j)$$

we want to get $R_{S^{(i)}}(\hat{w}_{S^{(i)}}) - R_{S^{(i)}}(\hat{w}_S)$
into the expression, just need to add and
subtract the appropriate terms

$$= \frac{1}{N}\left(\ell(\hat{w}_{S^{(i)}}, z_i) - \ell(\hat{w}_S, z_i)\right) \quad \underset{\text{Lipschitzness}}{\overset{\text{bound by}}{\leftarrow}}$$

$$+ \frac{1}{N}\left(\ell(\hat{w}_S, z_i') - \ell(\hat{w}_{S^{(i)}}, z_i')\right) \nwarrow$$

$$+ \left(\underbrace{R_{S^{(i)}}(\hat{w}_{S^{(i)}}) - R_{S^{(i)}}(\hat{w}_S)}_{\leq 0 \text{ by the definition of ERM}}\right)$$

Now putting it all together, we have

$$\leq \frac{2L}{N}\|\hat{w}_{S^{(i)}} - \hat{w}_S\|_2$$

$$\Rightarrow \frac{\alpha}{2}\|\hat{w}_{S^{(i)}} - \hat{w}_S\|_2^2 \leq \frac{2L}{N}\|\hat{w}_{S^{(i)}} - \hat{w}_S\|_2$$

$$\Rightarrow \|\hat{w}_{S^{(i)}} - \hat{w}_S\|_2 \leq \frac{4L}{\alpha N}$$

using Lipschitzness again

$$|\ell(\hat{w}_{S^{(i)}}, z_i') - \ell(\hat{w}_S, z_i')| \leq L \|\hat{w}_{S^{(i)}} - \hat{w}_S\|_2$$

$$\leq \frac{4L^2}{\alpha N}$$

Since the l.h.s. holds for any $S, S'$ and $i$, we have just shown the desired bound on uniform stability. ◩

To summarize, what we've shown is

strongly convex loss $\Rightarrow$ uniform stability $\Rightarrow$ average stability $\approx$ general. error

There is still one more part to the story

## Regularization

what can you do if your loss is convex and Lipschitz, but not strongly convex?

"Adding a strongly convex regularizer to the objective still approximately minimize the original objective."

Let $F_S(w) \overset{\Delta}{=} R_S(w) + \frac{\alpha}{2}\|w\|_2^2$, similarly
for $F(w) \overset{\Delta}{=} R(w) + \frac{\alpha}{2}\|w\|_2^2$

This is called $\ell_2$-regularization, weight
decay or Tikhonov regularization, etc

Moreover let

$$\hat{w}_S(\alpha) \overset{\Delta}{=} \underset{w \in \Omega}{\operatorname{argmin}} \; F_S(w)$$

be the ERM of the regularized objective

<u>Theorem 4</u> [Shalev-Shwartz et al] Suppose
$\ell$ is convex and L-Lipschitz. And let

$$S = (z_1, \ldots, z_N)$$

be an i.i.d. sample from D. Then

$$R(\hat{w}_S(\alpha)) - \underset{\underset{\text{minimizer of } R}{\uparrow}}{R^*} \leq 4\sqrt{\frac{L^2 B^2}{\delta N}}\left(1 + \frac{8}{\delta N}\right)$$

with probability $\geq 1-\delta$, where B is an

upper bound on the radius of $\Omega$ and we set

$$\alpha = \sqrt{\frac{16 L^2}{8 B^2 N}}$$

Proof: First notice that

$$f(w, z) \triangleq \ell(w, z) + \frac{\alpha}{2} \|w\|_2^2$$

is $\alpha$-strongly convex and $L + \alpha B$-Lipschitz

Now using Theorem 3 we have

$$F(\hat{w}_s(\alpha)) - F_s(\hat{w}_s(\alpha)) \leq \frac{4 (L + \alpha B)^2}{\alpha N} \quad (1)$$

We claim that

minimizer of F
↓

$$\mathbb{E}_s \left[ F_s(\hat{w}_s(\alpha)) \right] \leq F(w^*) \quad (2)$$

which is true because

$$F_s(\hat{w}_s(\alpha)) = \min_w F_s(w) \leq F_s(w^*)$$

and taking expectations on both sides

$$\mathbb{E}_S[F_S(\hat{w}_S(\alpha))] \leq \mathbb{E}_S[F_S(w^*)] = F(w^*)$$

Now combining (1) and (2) we have

$$\mathbb{E}_S[F(\hat{w}_S(\alpha)) - F(w^*)] \leq \frac{4(L + \alpha B)^2}{\alpha N}$$

Since the r.v. inside the expectation is nonnegative, we can apply Markov's bound

$$F(\hat{w}_S(\alpha)) - F(w^*) \leq \frac{4(L + \alpha B)^2}{\delta \alpha N}$$

with probability $\geq 1 - \delta$

Furthermore

$$F(w^*) \leq \inf_{w \in \Omega} R(w) + \frac{\alpha B^2}{2}$$

And so we have

$$R(\hat{w}_S(\alpha)) \leq F(\hat{w}_S(\alpha))$$

$$\leq R^* + \frac{\alpha B^2}{2} + \frac{4(L + \alpha B)^2}{8 \alpha N}$$

Now plugging in our choice of $\alpha$ completes the proof. ∎

## Implicit Regularization, Take 2

It turns out that explicit regularization is sometimes unneccessary

Theorem 5 [Hardt, Recht, Singer] Suppose $\ell(w, z)$ is convex, $B$-smooth and $L$-Lipschitz for every $z$. Then if we run SGD with step sizes

$$\eta_t \leq \frac{2}{B}$$

for $T$ timesteps, we have

$$\Delta_{unif}(SGD) \leq \frac{2L^2}{N} \sum_{t=1}^{T} \eta_t$$

In particular if we set $\eta_t \sim \frac{1}{\sqrt{N}}$ and $T = O(N)$ we get

$$\Delta_{unif}(SGD) \sim \frac{1}{\sqrt{N}} = \text{sampling error}$$

**Once again we get some of the benefits of regularization without needing to do it explicitly**

Note: SGD is a randomized algorithm so we need to modify the definition of uniform stability to

$$\sup_{s, s', i} \left| \mathop{\mathbb{E}}_A \left[ \ell(A(s), z_i') - \ell(A(s^{(i)}), z_i') \right] \right|$$

Intuition: At each step, SGD run on $S$ and $S^{(i)}$ has a $1 - \frac{1}{N}$ chance of picking the same sample, and if not can bound how much diverge

# Connections to Privacy

There are strong parallels between uniform stability and <u>differential privacy</u>:

> "The output of an algorithm should not change much if change one datapoint."

Operationally, this means that you can't tell whether a particular sample belongs to the training set based on the output of the learner

In a seminal work, Dwork, McSherry, Nissim and Smith introduced this concept:

<u>def:</u> [Dwork et al] An algorithm $\mathcal{A}$

is $(\varepsilon, \delta)$ - differentially private if for any set $U \subseteq \Omega$ we have

$$\mathbb{P}[A(S) \in U] \leq e^{\varepsilon} \mathbb{P}[A(S^{(i)}) \in U] + \delta$$

Now suppose you want to distinguish between two hypotheses

$H_0$: model was learned on $S$

$H_1$: model was learned on $S^{(i)}$

Then we can define

"false alarm" $\quad P_{FA} \triangleq$ Probability you rejected $H_0$ but it came from $H_0$

"missed detection" $\quad P_{MD} \triangleq$ Probability you retained $H_0$, it came from $H_1$

Theorem [Oh, Viswanath] If $A$ is
$(\varepsilon, \delta)$-differentially private

$$P_{FA} + e^{\varepsilon} P_{MD} \geq 1-\delta$$

$$e^{\varepsilon} P_{FA} + P_{MD} \geq 1-\delta$$

Thus when $\varepsilon$ and $\delta$ are small, you cannot beat random guessing by much

Notice there is a subtle difference:

### differential privacy

The parameters you learn do not depend much on any one data point

### uniform stability

The loss on any point $z$ does not change much, if you swap one data point

In particular, the nuissance params.

might not affect stability or generalization error, but might violate privacy!

There is a natural remedy:

"Inject noise into SGD"

Song, chaudhuri and Sarwate showed

**Theorem** [Song et al] SGD when you inject a large enough amount of noise, is differentiably private

**Intuition:** You can couple the trajectories when you run SGD on $S$ and $S^{(i)}$

Unfortunately:

reasonable level of privacy $\Rightarrow$ large amount of noise injected $\Rightarrow$ much worse learning performance

# Double Descent

For square loss, we can decompose our prediction error into

bias $\cong$ are our predicted labels off in expectation?

variance $\cong$ are we overfitting the data?

In particular, if we want to predict the responses

$$y = f(x) + \varepsilon \quad \leftarrow \substack{\text{mean zero,} \\ \text{variance } \sigma^2}$$

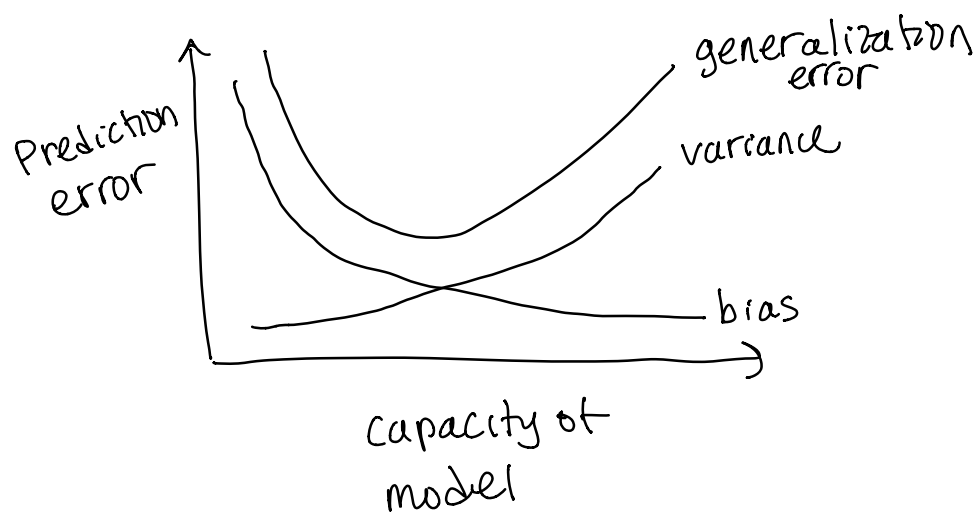and we learn a predictor $\hat{y} = \hat{f}(x; s)$ learned from the sample S

Then the square error can be written:

$$(y - \hat{y})^2 \quad = \quad \underbrace{\left( \mathbb{E}_S[\hat{f}(x; S)] - f(x) \right)^2}_{bias^2}$$

$$\mathbb{E}_S\left[ \underbrace{\left( \mathbb{E}_{S'}[\hat{f}(x; S')] - \hat{f}(x; S) \right)^2}_{variance} \right] + \underset{noise}{\sigma^2}$$
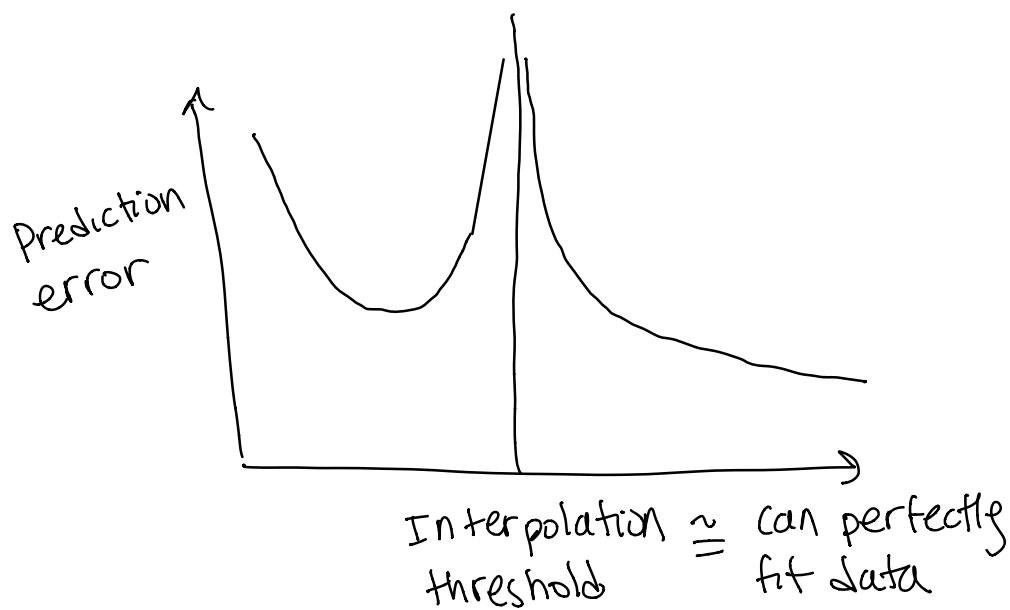
This is called the bias-variance decomposition

Classic generalization bounds suggest the following picture



But maybe things are more subtle...

An influential work of Belkin, Hsu, Ma
and Mandal showed the generalization
error curve is not always U-shaped

Instead they found



Prediction
error

Interpolation $\simeq$ can perfectly
threshold         fit data

In particular, you are better off
going beyond the point where the
training error is zero

Hastie, Montanari, Rosset and Tibshirani gave a rigorous treatment of this phenomenon in the context of regression

The standard model is

$$y_i = x_i^T \beta + \varepsilon_i \leftarrow \text{mean zero, variance } \sigma^2$$

$\uparrow$ mean zero, covariance $\Sigma$

Let $p$ = # of dimensions $\left.\right\}$ usual stat. notation
$n$ = # of data points

Let $\hat{\beta} = \arg\min \left\{ \|b\|_2 \mid b \text{ minimizes } \|y - Xb\|_2^2 \right\}$
$\uparrow$

Recall this is what gradient descent reaches when we initialize close to zero

Let's specialize to $\Sigma = I$

__Theorem__ [Hastie et al] Let $\Sigma = I$

and $\frac{P}{n} \to \gamma \in (0, \infty)$. Then

$$R(\hat{\beta}, \beta) \to \begin{cases} \sigma^2 \frac{\gamma}{1-\gamma} & \text{for } \gamma < 1 \\ r^2(1-\frac{1}{\gamma}) + \frac{\sigma^2}{\gamma-1} & \text{for } \gamma > 1 \end{cases}$$

where $r^2 = \|\beta\|_2^2$ and

$$R(\hat{\beta}, \beta) = \mathbb{E}\left[(x^\top \hat{\beta} - x^\top \beta)^2\right]$$

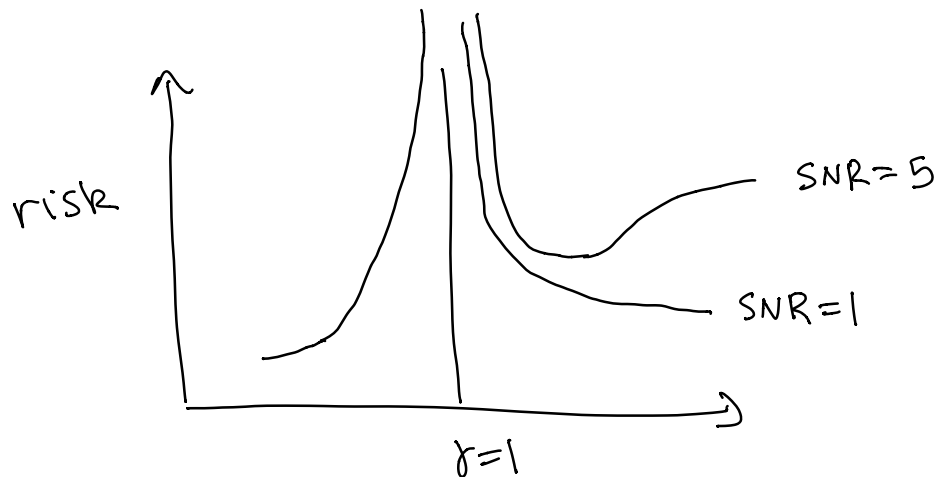__Interpretation__: Let $SNR \stackrel{\Delta}{=} \frac{r^2}{\sigma^2}$

First, we can always choose $\hat{\beta} = 0$, in which case

$$R(0, \beta) = r^2 \stackrel{\sim}{=} \text{null risk}$$

__Observation 1__: Regardless of SNR, as $\gamma \to \infty$ the risk converges to the null risk

## Observation 2: If SNR < 1, the risk is always a decreasing function of $\gamma$ on $(1, \infty)$

In a picture

J