<u>Reminder:</u>

$$\hat{R}_S[h] = \frac{1}{n} \sum_{i=1}^{n} \ell(h(x_i), y_i) \leftarrow \text{empirical risk}$$

$$R[h] = E_{z \sim D} \ell(h(x), y) \leftarrow \text{true risk}$$

## Last time:

Empirical risk concentration.

- <u>Main message:</u> Empirical risk is within $\varepsilon$ from true if

$$\#\text{Samples} \geqslant O\left(\frac{\#\text{params}}{\varepsilon^2}\right)$$

However, There are very sophisticated techniques to extend this beyond finite classes. eg. • VC-dimension

  • Rademacher complexity.

In future lectures we will see concentration bounds that are "algorithm-specific"

## Lecture 3:

- From statistical bounds to optimization
- Computational aspects of the ERM
- Examples of loss fns

What do concentration bounds tell us?

- Lets assume that

$$\hat{R}_S[h] \leq R[h] + \varepsilon, \quad \forall h \in \mathcal{H}$$

with prob. $1-\delta$

But we care about a "special" $h$:

- $\hat{h}^* = \arg\min_{h \in \mathcal{H}} \hat{R}[h] \quad \leftarrow \text{ERM}$

and its true performance

$$R[\hat{h}^*] = E_z \, \ell(\hat{h}^*(x), y)$$

Then, if we have concentration $\forall h \in \mathcal{H}$

$$\Rightarrow \quad R[\hat{h}^*] = \hat{R}[\hat{h}^*] + \left( \hat{R}[\hat{h}^*] - R[\hat{h}^*] \right)$$

$$= \hat{R}[\hat{h}^*] + \left( R[h^*] - \hat{R}[h^*] \right)$$

$$\leq \hat{R}[h^*] + \varepsilon \quad \text{w.p. } 1-\delta$$

If the ER concentrates Then

$$R[\hat{h}^*] \leq \hat{R}[\hat{h}^*] + \varepsilon$$

But also we can relate $R[\hat{h}^*]$ to the best predictor in $\mathcal{H}$

$$h^* = \underset{h \in \mathcal{H}}{\arg\min} \; R[h]$$

We have that

$$R[\hat{h}^*] \leq \hat{R}[\hat{h}^*] + \varepsilon$$
$$\leq \hat{R}[h^*] + \varepsilon$$
$$\leq R[h^*] + \underbrace{\hat{R}[h^*] - R[h^*]}_{\varepsilon} + \varepsilon$$
$$\leq R[h^*] + 2\varepsilon$$

Hence, we can argue about the best possible predictor via the performance of the ERM, assuming concentration

The above are a brief preview
of "why ERM is a good idea".

But what does it look like?

Examples:

- Regression:

  - linear: $\min_w \| Xw - y \|^2$ $\begin{bmatrix} +\lambda\|w\|_2^2 \text{ ridge} \\ +\lambda\|w\|_1 \text{ lasso} \end{bmatrix}$

    $\equiv \min_w \frac{1}{n} \sum (x_i^T w - y_i)^2$

  - nonlinear (e.g. nn)

    $\min \frac{1}{n} \sum (y_i^2 - h(x_i; w))^2$

    $h(x_i; W) = 6(W_L \, 6(W_{L-1} \, 6(\dots \, 6(W_1 x)))$

- Classification:
  - Binary: $\min \frac{1}{n} \sum_{i=1}^{n} \log(1 + e^{-y_i x_i^T w})$

    $\min \frac{1}{n} \sum_{i=1}^{n} \max\{0, 1 - y_i x_i^T w\}$

- **Multiclass:** (for one sample)

$$-\sum_{c=1}^{M} y_{i,c} \log\left([h(w; x_i)]_c\right)$$

# Back to optimization:

We want to solve

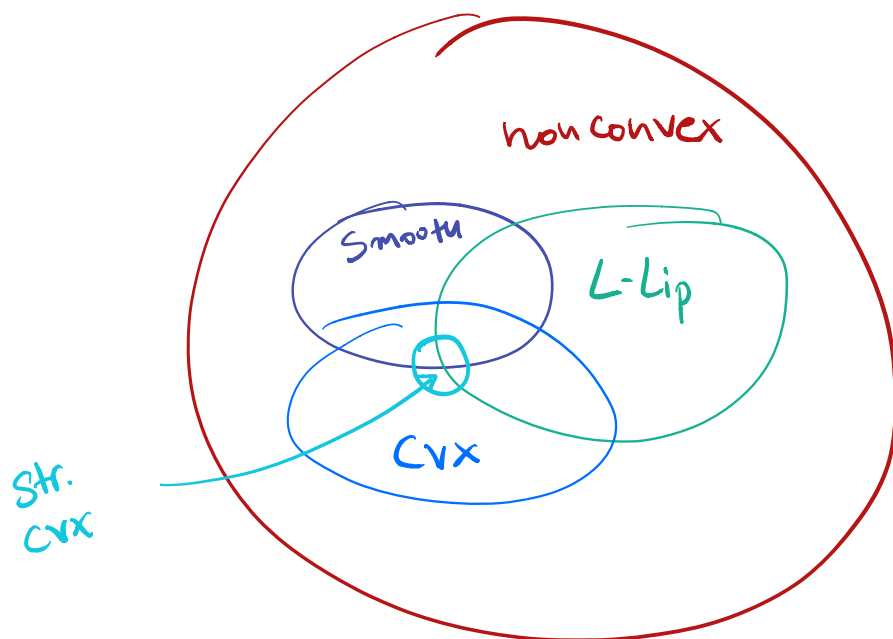$$\min_{w} \frac{1}{n} \sum_{i=1}^{n} l_i(w) + \lambda \cdot R(w)$$

- Q: When can we solve it?
- Q: How fast?

Remark: The more you know about the structure of the problem the more we can say about "solvability" and "scalability".

Informal Theorem: In the general case, ERM is NP-Hard.
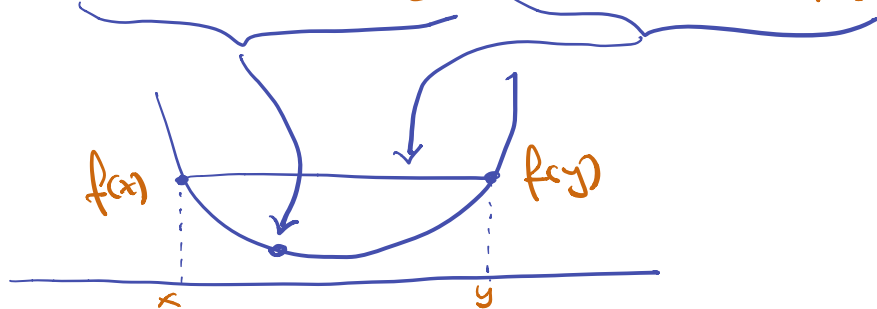
let's put some structure

Families of functions



nonconvex
Smooth
L-Lip
Cvx
Str. Cvx

• Convex:

$$f(a \cdot \vec{x} + (1-a)\vec{y}) \leq a \cdot f(\vec{x}) + (1-a) f(\vec{y})$$

also

$f(x) \geq f(y) + \langle \nabla f(y), x-y \rangle$

$f(x)$    $f(y)$

$x$    $y$

Convexity makes our lives easy, eg

$\min f(\omega)$  solvable in poly-time

- **Important property of cvx functions:**

  Every <u>local min</u> ≡ <u>global min</u>

<u>L - Lipschitz:</u> ("fns that don't change fast")

$$\forall \ x, y \qquad |f(x) - f(y)| \leq L\|x-y\|$$

<u>B-smooth:</u> ("fns with grads that don't change fast")

$$\forall \ x, y \qquad \|\nabla f(x) - \nabla f(y)\| \leq B \cdot \|x-y\|$$

<u>strongly convex:</u> ("the best kind of fns")

$$f(x) - f(y) \leq \langle \nabla f(x), x-y \rangle - \frac{\partial}{2}\|x-y\|^2$$

# Examples:

- **Convex:**
  - $\|x\|^2, \|x\|, \log(1+\exp(x)), \max\{0, 1-x\}$
  - if $g(\cdot)$ is cvx, then $g(\vec{x}^T w + b)$ is

  eg: $\log(1 + \exp(-y \langle w, x \rangle))$
  $(w^T x - b)^2 \dots$

  - $\max_i f_i(x)$    (if $f_i$ are cvx)

  - $\sum_i w_i f_i(x)$    (if $f_i$ are cvx)

  - $\dots$

- **L-Lip:**
  - $|x|$ is $1$-Lip.
  $f(x) = \log(1 + \exp(x))$ is $1$-Lip

  - $x^2$ <u>is not</u> Lipschitz
  unless $|x| \leq p$ when it is $p$-Lip

  - $f(w) = x^T w + b$ is $\|x\|$-Lip.

  - $f(x) = g_1(g_2(x))$. If $g_1$ is $L_1$-Lip.
  $g_2$ is $L_2$-Lip
  $\Rightarrow f$ is $L_1 \cdot L_2$-Lip.

- $g(x^T w + b) \implies \|x\| \cdot L_g - \text{Lip.}$

- If $\|\nabla f(w)\| \le L \implies f$ is $L$-Lip.

## Smooth:

- $|x|^2$ is $2$-smooth

- $\log(1 + e^x)$ is $\frac{1}{4}$-smooth

- If $g$ is $\beta$-smooth
  $f(w) = g(w^T x + b)$ is $\beta \|x\|^2$-smooth

- $f(w) = \log(1 + e^{-y\langle w, x\rangle})$  $\frac{\|x\|^2}{4}$-smooth

## Strongly Convex:

- $f$ is $\lambda$-str convex if
  $f(w) - \frac{\lambda}{2}\|w\|^2$ is convex

- eg. $\sum_{i=1}^{n} \log(1 + e^{-y\langle w, x\rangle}) + \frac{\lambda}{2}\|w\|^2$

  $\vdots$

## Next time:

- why is convexity useful?

- How to exploit it algorithmically?

- Gradient Methods